



CS-463 Information Retrieval Systems

Μοντέλα Ανάκτησης (Retrieval Models)

Part C

Yannis Tzitzikas

University of Crete

CS-463, Spring 05

Lecture : 5

Date : 8-3-2005



Διάρθρωση Διάλεξης

PART (A)

- Ανάκτηση και Φιλτράρισμα
- Εισαγωγή στα Μοντέλα Αντήλησης
- Κατηγορίες Μοντέλων
- Exact vs Best Match
- Τα κλασσικά μοντέλα ανάκτησης
 - Το Boolean Μοντέλο
 - Στατιστικά Μοντέλα - Βάρυνση Όρων
 - Το Διανυσματικό Μοντέλο
 - Το Πιθανοκρατικό Μοντέλο

PART (B): **Εναλλακτικά μοντέλα**

- (I) Συνολοθεωρητικά μοντέλα
 - Fuzzy Retrieval Model
 - Extended Boolean Model
- (II) Αλγεβρικά Μοντέλα
 - Latent Semantic Indexing
 - Neural Network Model

PART (C):

- (III) Πιθανοκρατικά Μοντέλα
 - Bayesian Network Model
 - Inference Network Model



Εναλλακτικά Πιθανοκρατικά Μοντέλα

- Probability Theory
 - Semantically clear
 - Computationally clumsy
- Why Bayesian (Belief) Networks?
 - Clean formalism to combine distinct sources of evidence
 - past queries, past feedback cycles, distinct query formulations
 - Modularize the world (dependencies)
- Bayesian Network Models for IR:
 - Inference Network (Turtle & Croft, 1991)
 - Belief Network (Ribeiro-Neto & Muntz, 1996)



Bayesian Inference

Basic Probability Axioms:

$$0 \leq P(A) \leq 1 ;$$

$$P(\text{sure})=1;$$

$$P(A \vee B)=P(A)+P(B) \quad \text{if } A \text{ and } B \text{ are mutually exclusive}$$

Other formulations

- $P(A)=P(A \wedge B)+P(A \wedge \neg B)$

- $P(A)=\sum_{B_i} P(A \wedge B_i)$, where B_{i,v_i} is a set of exhaustive and mutually exclusive events

- $P(A) + P(\neg A) = 1$

- $P(A|K)$ belief in A given the knowledge K

- if $P(A|B)=P(A)$, we say: A and B are *independent*

- if $P(A|B \wedge C)=P(A|C)$, we say: A and B are conditionally independent, given C

- $P(A \wedge B)=P(A|B)P(B)$

- $P(A)=\sum_{B_i} P(A | B_i)P(B_i)$



Bayesian Inference

Bayes' Rule : the heart of Bayesian techniques

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

where,

- H : a hypothesis
- e : is an evidence
- P(H) : prior probability
- P(H|e) : posterior probability
- P(e|H) : probability of e if H is true
- P(e) : a normalizing constant, then we write: $P(H|e) \sim P(e|H)P(H)$



Bayesian Networks

Bayesian networks are **directed acyclic graphs (DAGS)** in which the **nodes** represent *random variables*, the **arcs** portray *causal relationships* between these variables, and the **strengths** of these causal influences are expressed by *conditional probabilities*.

y_i : parent nodes (in this case, root nodes)

x : child node

y_i cause x

Y the set of parents of x

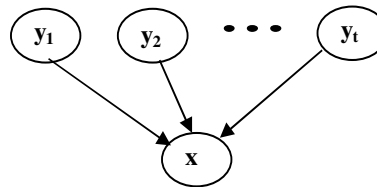
The influence of Y on x

can be quantified by any function

$F(x,Y)$ such that $\sum_{\forall x} F(x,Y) = 1$

$$0 \leq F(x,Y) \leq 1$$

For example, $F(x,Y)=P(x|Y)$





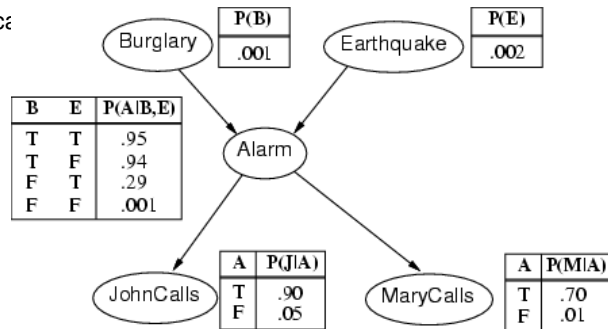
Παράδειγμα

I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?

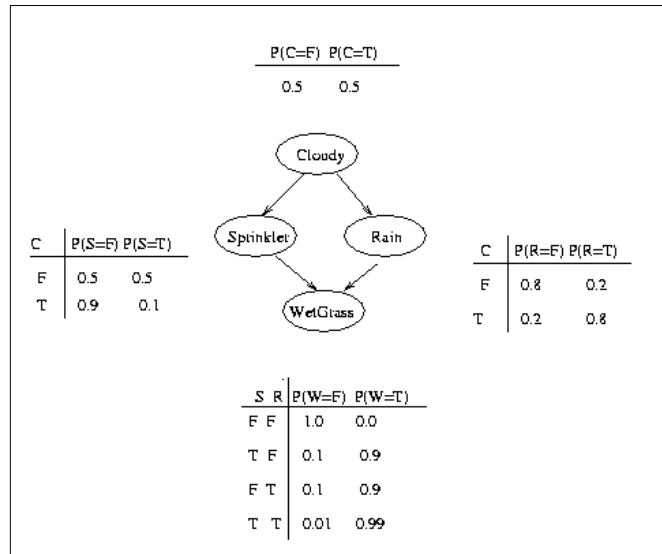
Variables: *Burglary, Earthquake, Alarm, JohnCalls, MaryCalls*

Network topology reflects "causal" knowledge:

- A burglar can cause the alarm
- An earthquake can set the alarm
- The alarm can cause Mary to call
- The alarm can cause John to call



Παράδειγμα



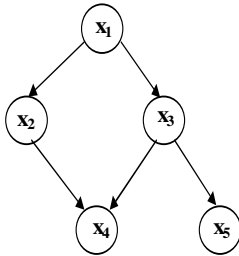


Bayesian Networks and Joint Probabilities

The full joint distribution is defined as the product of the local conditional distributions:

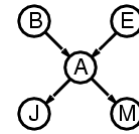
$$P(x_1, \dots, x_n) = \prod_{i=1, n} P(x_i | \text{Parents}(x_i)) \square$$

Examples:



$$P(x_1, x_2, x_3, x_4, x_5) = P(x_1) P(x_2 | x_1) P(x_3 | x_1) P(x_4 | x_2, x_3) P(x_5 | x_3)$$

$P(x_1)$: prior probability of the root node



$$P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) = P(j | a) P(m | a) P(a | \neg b, \neg e) P(\neg b) P(\neg e)$$

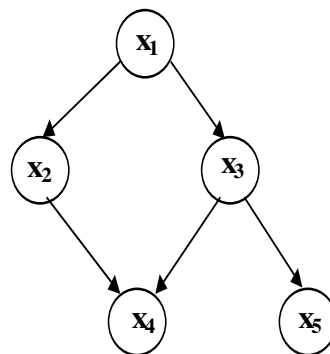


Also

In a Bayesian network each variable x is conditionally independent of all its non-descendants, given its parents.
(ανεξάρτητης των μη απογόνων δοθέντων των πατέρων)

For example:

$$P(x_4, x_5 | x_2, x_3) = P(x_4 | x_2, x_3) P(x_5 | x_3)$$





Information Retrieval Models

Inference Network Model



Inference Network Model

- Επιστημολογική άποψη του προβλήματος της ΑΠ
- Σχολές σκέψης στις πιθανότητες:
 - *freqüentist*:
πιθανότητα = στατιστικό μέγεθος σχετιζόμενο με την τύχη
 - *epistemological*:
πιθανότητα = βαθμός πίστης του οποίου η περιγραφή μπορεί να είναι απαλλαγμένη από στατιστικά πειράματα
- Η προσέγγιση του Inference Network Model:
 - Επιστημολογική χρήση των πιθανοτήτων



Inference Network Model: Τυχαίες μεταβλητές

- **Εγγράφων**
 - Η τυχαία μεταβλητή που σχετίζεται με ένα έγγραφο d_j παριστάνει το **συμβάν παρατήρησης** του εγγράφου αυτού
- **Όρων**
 - Η παρατήρηση ενός εγγράφου είναι η αιτία για **αυξημένη πίστη** στις τυχαίες μεταβλητές που αντιστοιχούν στους όρους που περιέχει το κείμενο
- **Επερωτήσεων**
 - εκφράζει το βαθμό ικανοποίησης της επερώτησης
- Όλες οι τυχαίες μεταβλητές (d, k, q) είναι **δυναδικές** (0 ή 1) !

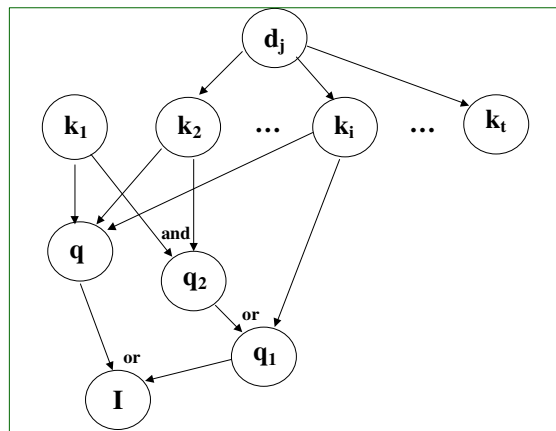


Inference Network Model: Δομή (documents => terms => queries)

d_j has index terms $k_2, k_i,$ and k_t
 q has index terms $k_1, k_2,$ and k_i
 q_1 and q_2 model boolean formulation
 $q_1 = ((k_1 \wedge k_2) \vee k_i);$
 $I = (q \vee q_1)$

Nodes

documents (d_j)
index terms (k_i)
queries ($q, q_1,$ and q_2)
user information need (I)



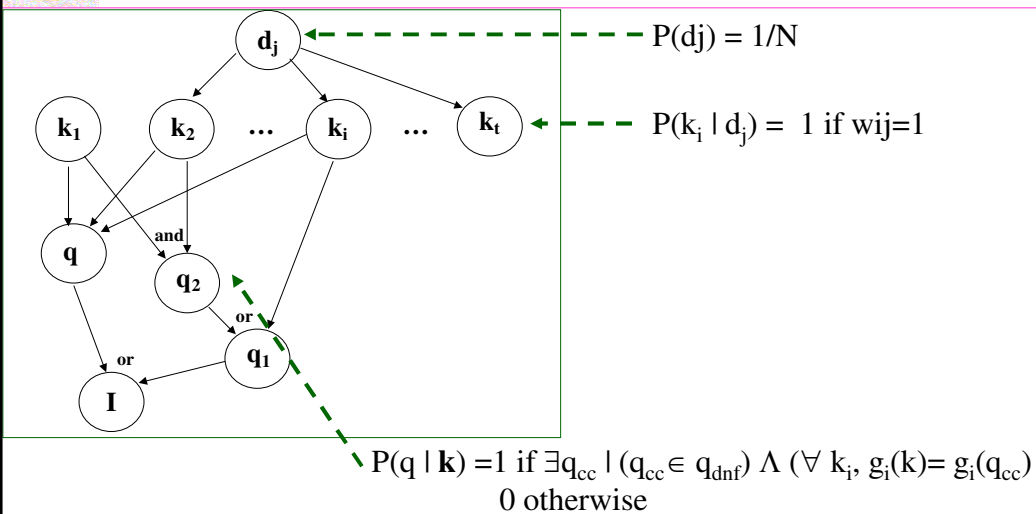


Inference Network Model: Ranking

- Όλες οι τυχαίες μεταβλητές (d, k, q) είναι **δυναδικές** (0 ή 1) !
- **Κατάταξη** $\text{Rank}(q, d_j)$
 - $\text{Rank}(q, d_j) = P(q \wedge d_j)$
 - q and d_j είναι συντομογραφίες του $q=1$ and $d_j=1$
 - $\text{Rank}(q, d_j)$:
 - expresses how much evidential support the observation of d_j provides to the query q
 - $\text{Rank}(q, d_j) = P(q=1 \wedge d_j=1)$
 - (d_j stands for a state where $d_j = 1$ and $\forall i \neq j \Rightarrow d_i = 0$, because we observe one document at a time)



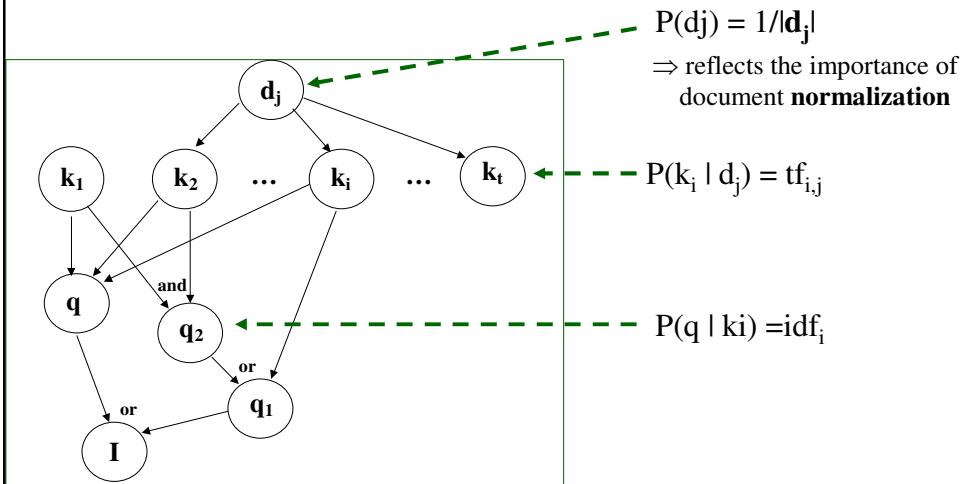
How an inference network can be tuned to subsume the Boolean Model



\Rightarrow one of the conjunctive components of the query **must be** matched by the active index terms in \mathbf{k}



Inference network and TF-IDF



$$P(d_j) = 1/|d_j|$$

⇒ reflects the importance of document **normalization**

$$P(k_i | d_j) = tf_{i,j}$$

$$P(q | k_i) = idf_i$$

⇒ however the obtained ranking is distinct of the one provided by the vector model



Inference Network Model: Σύνοψη

- Επιστημολογική προσέγγιση
- Τοπολογία δικτύου $d \rightarrow k \rightarrow q$
- Εκφραστική ικανότητα
 - Συλλαμβάνει το Boolean Model (με δυαδικές τυχαίες μεταβλητές)
 - Μπορεί να προσεγγίσει αρκετά το διανυσματικό
 - Δυνατότητα συνδυασμού πολλών αποδεικτικών πηγών
- Συστήματα βασισμένα σε αυτό το μοντέλο
 - Inquiry system



Information Retrieval Models

Belief Network Model



Belief Network Model

- Όπως και στο Inference Network Model
 - Επιστημολογική άποψη του προβλήματος της ΑΠ
 - Τυχαίες μεταβλητές για έγγραφα, όρους και επερωτήσεις
- Αντίθετα με το Inference Network Model, εδώ έχουμε
 - πιο ξεκάθαρο δειγματικό χώρο
 - συνολοθεωρητική προσέγγιση
 - διαφορετική τοπολογία δικτύου



Belief Network Model: The Probability Space

Δειγματικός Χώρος: $\mathbf{K}=\{k_1, k_2, \dots, k_t\}$ (sample space)

Κάθε όρος k_i είναι μια στοιχειώδης έννοια (elementary concept)

Κάθε σύνολο $u \subset \mathbf{K}$ είναι μια έννοια

Σε κάθε $u \subset \mathbf{K}$ αντιστοιχίζεται ένα δυαδικό διάνυσμα $\mathbf{k}=(k_1, k_2, \dots, k_t)$
τ.ω. $w_i=1 \Leftrightarrow k_i \in u$

k_i a binary random variable associated with the index term k_i , ($k_i = 1$
 $\Leftrightarrow w_i=1 \Leftrightarrow k_i \in u$)



Belief Network Model: Documents and Queries

Εγγραφο d :

ένα υποσύνολο του \mathbf{K} (που περιέχει τους όρους του d)

Επερώτηση q :

ένα υποσύνολο του \mathbf{K} (που περιέχει τους όρους της q)

Ορίζουμε μια κατανομή πιθανότητας P στο \mathbf{K} ως εξής:

$P(c)=\sum_u P(c|u) P(u)$ // the *degree of coverage* of the space \mathbf{K} by c
 c : a generic concept representing a document or a query

Αρχικά (που δεν γνωρίζουμε τα $P(u)$), υποθέτουμε ότι:

$P(u)=(1/2)^t = 1/2^{|\mathbf{K}|}$

Belief Network Model:
Τοπολογία Δικτύου και Σημασιολογία

$P(q)=1$ means that q covers completely the concept space

$P(d_j)=1$ means that d_j covers completely the concept space

So user queries and documents are modeled as **subsets** of index terms.

The concept space K works as the common sample space.

CS-463, Information Retrieval
 Yannis Tzitzikas, U. of Crete, Spring 2005

137

Belief Network Model:
Ranking : $P(d|q)$

$P(d | q)$
 Εκφράζει το βαθμό κάλυψης του d από την επερώτηση q

Όσο πιο μεγάλος είναι, τόσο πιο συναφές είναι το d με το q

CS-463, Information Retrieval
 Yannis Tzitzikas, U. of Crete, Spring 2005

138

**Belief Network Model:
Computing $P(d|q)$**

$P(d | q)$
Εκφράζει το βαθμό κάλυψης του d από την επερώτηση q

$P(d_j|q) = P(d_j \wedge q) / P(q)$ (εξ' ορισμού: θεώρημα Bayes)
 $\sim P(d_j \wedge q)$ (διότι ο παρανομαστής $P(q)$ είναι ίδιος για όλα τα d)
 $\sim \sum_u P(d_j \wedge q | u) P(u)$ (εξ' ορισμού)
 $\sim \sum_u P(d_j | u) P(q | u) P(u)$ (Bayesian networks)
 $\sim \sum_k P(d_j | \mathbf{k}) P(q | \mathbf{k}) P(\mathbf{k})$

Αρκεί να ορίσουμε αυτά

CS-463, Information Retrieval, Yannis Tzitzikas, U. of Crete, Spring 2005 139

Defining $P(q|k)$ and $P(d|k)$ to model the vector model

$P(q | k_i) = \frac{w_{iq}}{\sqrt{\sum_{i=1}^t w_{iq}^2}}$ if q contains k_i (otherwise = 0)

$P(d_j | k_i) = \frac{w_{ij}}{\sqrt{\sum_{i=1}^t w_{ij}^2}}$ if d_j contains k_i (otherwise = 0)

In this way we get the ranking of the classic Vector Space Model !

CS-463, Information Retrieval, Yannis Tzitzikas, U. of Crete, Spring 2005 140



Belief Network Model: Σύνοψη

- Επιστημολογική προσέγγιση
- Τοπολογία δικτύου $d \leftarrow k \rightarrow q$
 - (το inference network είχε δομή $d \rightarrow k \rightarrow q$)
- Εκφραστική ικανότητα
 - Συλλαμβάνει το Boolean Model (με δυαδικές τυχαίες μεταβλητές)
 - Συλλαμβάνει το Διανυσματικό Μοντέλο
 - Μπορεί αναπαράγει οποιαδήποτε διάταξη του Inference Network (το αντίστροφο δεν ισχύει)
 - Δυνατότητα συνδυασμού πολλών αποδεικτικών πηγών
- Inference Network Model:
 - το πρώτο και πιο γνωστό
 - επιτυχημένη χρήση στο σύστημα Inquiry



Διάρθρωση

PART (A)

- Ανάκτηση και Φιλτράρισμα
- Εισαγωγή στα Μοντέλα Αντιληψης
- Κατηγορίες Μοντέλων
- Exact vs Best Match
- Τα κλασσικά μοντέλα ανάκτησης
 - Το Boolean Μοντέλο
 - Στατιστικά Μοντέλα - Βάρυνση Όρων
 - Το Διανυσματικό Μοντέλο
 - Το Πιθανοκρατικό Μοντέλο

PART (B): Εναλλακτικά μοντέλα

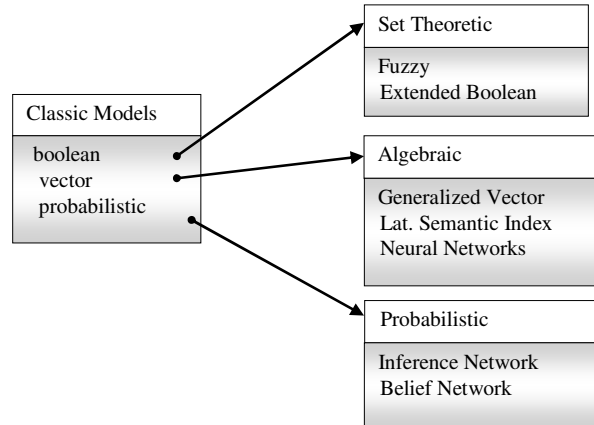
- (I) Συνολοθεωρητικά μοντέλα
 - Fuzzy Retrieval Model
 - Extended Boolean Model
- (II) Αλγεβρικά Μοντέλα
 - Latent Semantic Indexing
 - Neural Network Model

PART (C):

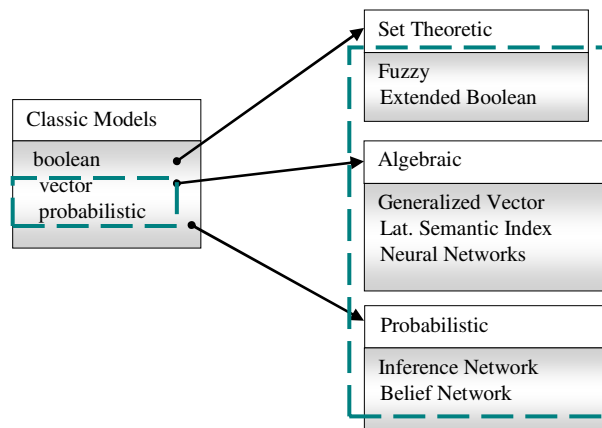
- (III) Πιθανοκρατικά Μοντέλα
 - Bayesian Network Model
 - Inference Network Model



Ταξινόμια Μοντέλων που εξετάσαμε



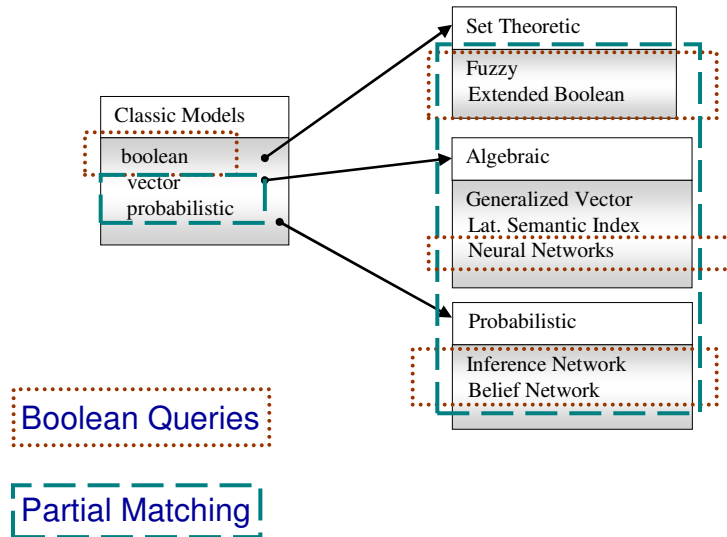
Ταξινόμια Μοντέλων που εξετάσαμε



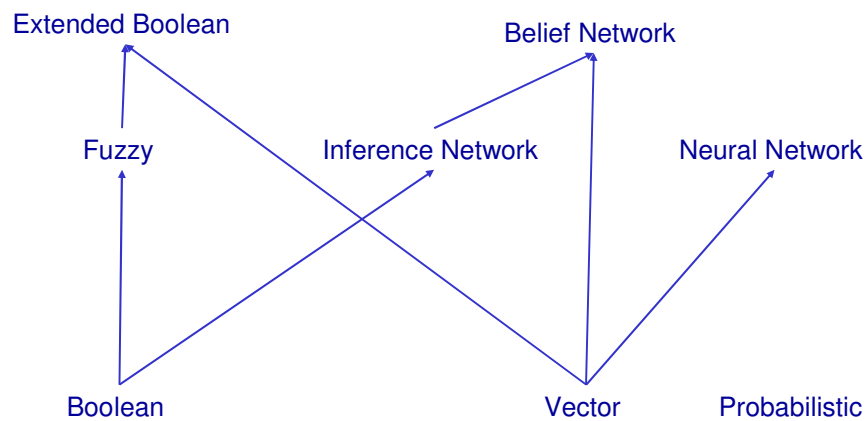
Partial Matching



Ταξινόμια Μοντέλων που εξετάσαμε



WRT Expressive Power (incomplete)





Άλλοι τύποι Μοντέλων Ανάκτησης που θα δούμε αργότερα



Αργότερα

- Μοντέλα Ανάκτησης **Δομημένων** Εγγράφων
 - Non Overlapping Lists
 - Proximal Nodes
 - Retrieval Models for XML
- Μοντέλα Ανάκτησης **Ιστοσελίδων**
 - Έμφαση στους συνδέσμους
- Μοντέλα Ανάκτησης **Πολυμέσων**
- Μοντέλα Βασισμένα στη **Λογική**
 - Carlo Meghini and Umberto Straccia, A Relevance Terminological Logic for Information Retrieval, Proceedings of SIGIR'96, Zurich, Switzerland, 1996