

Chapter 4: Advanced IR Models

4.1 Probabilistic IR

4.1.1 Principles

4.1.2 Probabilistic IR with Term Independence

4.1.3 Probabilistic IR with 2-Poisson Model (Okapi BM25)

4.1.4 Extensions of Probabilistic IR

4.2 Statistical Language Models

4.3 Latent-Concept Models

4.1.1 Probabilistic Retrieval: Principles

[Robertson and Sparck Jones 1976]

Goal:

Ranking based on $\text{sim}(\text{doc } d, \text{query } q) =$
 $P[R|d] = P [\text{doc } d \text{ is relevant for query } q \mid$
 $d \text{ has term vector } X_1, \dots, X_m]$

Assumptions:

- **Relevant and irrelevant documents differ in their terms.**
- **Binary Independence Retrieval (BIR) Model:**
 - **Probabilities for term occurrence are pairwise independent for different terms.**
 - **Term weights are binary $\in \{0,1\}$.**
- **For terms that do not occur in query q the probabilities for such a term occurring are the same for relevant and irrelevant documents.**

4.1.2 Probabilistic IR with Term Independence: Ranking Proportional to Relevance Odds

$$\begin{aligned} \text{sim}(d, q) &= O(R | d) = \frac{P[R | d]}{P[\neg R | d]} && \text{(odds for relevance)} \\ &= \frac{P[d | R] \times P[R]}{P[d | \neg R] \times P[\neg R]} && \text{(Bayes' theorem)} \\ &\sim \frac{P[d | R]}{P[d | \neg R]} = \prod_i \frac{P[X_i | R]}{P[X_i | \neg R]} && \text{(independence or linked dependence)} \\ \text{sim}(d, q)' &= \log \prod_{i \in q} \frac{P[X_i | R]}{P[X_i | \neg R]} && \text{(} X_i = 1 \text{ if } d \text{ includes } \\ & && \text{i-th term, 0 otherwise)} \\ &= \sum_{i \in q} \log P[X_i | R] - \log P[X_i | \neg R] \end{aligned}$$

Probabilistic Retrieval: Ranking Proportional to Relevance Odds (cont.)

$$= \sum_{i \in q} \log (p_i^{X_i} (1 - p_i)^{1 - X_i}) - \log (q_i^{X_i} (1 - q_i)^{1 - X_i}) \quad (\text{binary features})$$

with estimators $p_i = P[X_i = 1 | R]$ and $q_i = P[X_i = 1 | \neg R]$

$$= \sum_{i \in q} \log \left(\frac{p_i^{X_i} (1 - p_i)}{(1 - p_i)^{X_i}} \right) - \log \left(\frac{q_i^{X_i} (1 - q_i)}{(1 - q_i)^{X_i}} \right)$$

$$= \sum_{i \in q} X_i \log \frac{p_i}{1 - p_i} + \sum_{i \in q} X_i \log \frac{1 - q_i}{q_i} + \sum_{i \in q} \log \frac{1 - p_i}{1 - q_i}$$

$$\sim \sum_{i \in q} X_i \log \frac{p_i}{1 - p_i} + \sum_{i \in q} X_i \log \frac{1 - q_i}{q_i} = \text{sim}(d, q)''$$

Probabilistic Retrieval: Robertson / Sparck Jones Formula

Estimate p_i und q_i based on training sample
(query q on small sample of corpus) or based on
intellectual assessment of first round's result (*relevance feedback*):

Let N be #docs in sample,
 R be # relevant docs in sample
 n_i #docs in sample that contain term i ,
 r_i # relevant docs in sample that contain term i

$$\Rightarrow \text{Estimate: } p_i = \frac{r_i}{R} \quad q_i = \frac{n_i - r_i}{N - R}$$

$$\text{or: } p_i = \frac{r_i + 0.5}{R + 1} \quad q_i = \frac{n_i - r_i + 0.5}{N - R + 1} \quad (\text{Lidstone smoothing with } \lambda=0.5)$$

$$\Rightarrow \text{sim}(d, q)'' = \sum_i X_i \log \frac{r_i + 0.5}{R - r_i + 0.5} + \sum_i X_i \log \frac{N - n_i - R + r_i + 0.5}{n_i - r_i + 0.5}$$

$$\Rightarrow \text{Weight of term } i \text{ in doc } d: \log \frac{(r_i + 0.5) (N - n_i - R + r_i + 0.5)}{(R - r_i + 0.5) (n_i - r_i + 0.5)}$$

Probabilistic Retrieval: tf*idf Formula

Assumptions (without training sample or relevance feedback):

- p_i is the same for all i .
- Most documents are irrelevant.
- Each individual term i is infrequent.

This implies:

- $\sum_i X_i \log \frac{p_i}{1-p_i} = c \sum_i X_i$ with constant c
- $q_i = P[X_i = 1 | \neg R] \approx \frac{df_i}{N}$
- $\frac{1-q_i}{q_i} = \frac{N-df_i}{df_i} \approx \frac{N}{df_i}$

$$\begin{aligned} \Rightarrow \text{sim}(d, q) &= \sum_i X_i \log \frac{p_i}{1-p_i} + \sum_i X_i \log \frac{1-q_i}{q_i} \\ &\approx c \sum_i X_i + \sum_i X_i \text{idf}_i \end{aligned}$$

Scalar product over the product of tf and dampend idf values for query terms

Example for Probabilistic Retrieval

Documents with relevance feedback:

q: t1 t2 t3 t4 t5 t6

	t1	t2	t3	t4	t5	t6	R
d1	1	0	1	1	0	0	1
d2	1	1	0	1	1	0	1
d3	0	0	0	1	1	0	0
d4	0	0	1	0	0	0	0
ni	2	1	2	3	2	0	
ri	2	1	1	2	1	0	
pi	5/6	1/2	1/2	5/6	1/2	1/6	
qi	1/6	1/6	1/2	1/2	1/2	1/6	

} R=2, N=4

Score of new document d5 (with Lidstone smoothing with $\lambda=0.5$):

$$d5 \cap q: \langle 1 \ 1 \ 0 \ 0 \ 0 \ 1 \rangle \rightarrow \text{sim}(d5, q) = \log 5 + \log 1 + \log 0.2 \\ + \log 5 + \log 5 + \log 5$$

$$\text{sim}(d, q)'' = \sum_i X_i \log \frac{p_i}{1 - p_i} + \sum_i X_i \log \frac{1 - q_i}{q_i}$$

Laplace Smoothing (with Uniform Prior)

Probabilities p_i and q_i for term i are estimated

by MLE for binomial distribution

(repeated coin tosses for relevant docs, showing term i with p_i ,

Repeated coin tosses for irrelevant docs, showing term i with q_i)

To avoid overfitting to feedback/training,

the estimates should be smoothed

(e.g. with uniform prior):

Instead of estimating $p_i = k/n$ estimate (Laplace's law of succession):

$$p_i = (k+1) / (n+2)$$

or with heuristic generalization (Lidstone's law of succession):

$$p_i = (k+\lambda) / (n+2\lambda) \text{ with } \lambda > 0 \text{ (e.g. } \lambda=0.5)$$

And for multinomial distribution (n times w -faceted dice) estimate:

$$p_i = (k_i + 1) / (n + w)$$

4.1.3 Probabilistic IR with Poisson Model (Okapi BM25)

Generalize term weight $w = \log \frac{p(1-q)}{q(1-p)}$

into $w = \log \frac{p_{tf} q_0}{q_{tf} p_0}$

with p_j, q_j denoting prob. that term occurs j times in rel./irrel. doc

Postulate Poisson (or Poisson-mixture) distributions:

$$p_{tf} = e^{-\lambda} \frac{\lambda^{tf}}{tf!} \quad q_{tf} = e^{-\mu} \frac{\mu^{tf}}{tf!}$$

Okapi BM25

Approximation of Poisson model by similarly-shaped function:

$$w := \log \frac{p(1-q)}{q(1-p)} \cdot \frac{tf}{k_1 + tf}$$

finally leads to Okapi BM25 (which achieved best TREC results):

$$w_j(d) := \frac{(k_1 + 1)tf_j}{k_1 \left((1-b) + b \frac{\text{length}(d)}{\text{avgdoclength}} \right) + tf_j} \cdot \log \frac{N - df_j + 0.5}{df_j + 0.5}$$

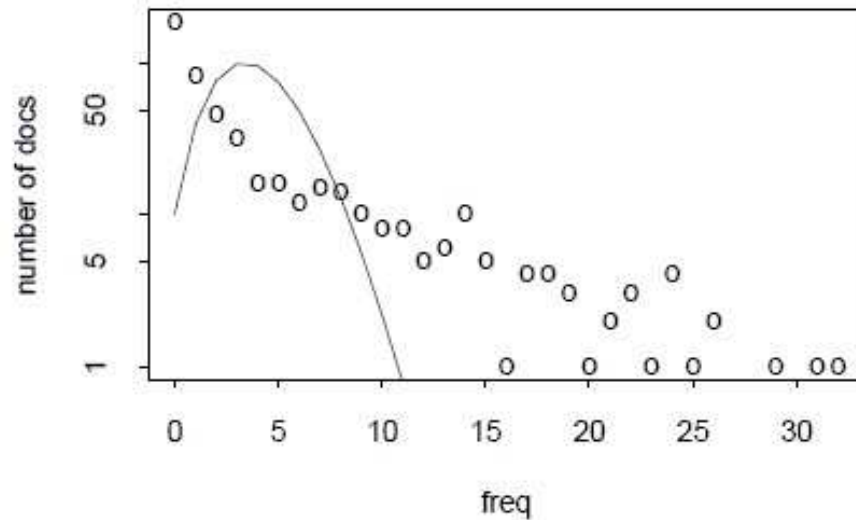
or in the most comprehensive, tunable form:

$$\text{score}(d, q) := \sum_{j=1..|q|} \log \frac{N - df_j + 0.5}{df_j + 0.5} \cdot \frac{(k_1 + 1)tf_j}{k_1 \left((1-b) + b \frac{\text{len}(d)}{\Delta} \right) + tf_j} \cdot \frac{(k_3 + 1)qtf_j}{k_3 + tf_j} + k_2 |q| \frac{\Delta - \text{len}(d)}{\Delta + \text{len}(d)}$$

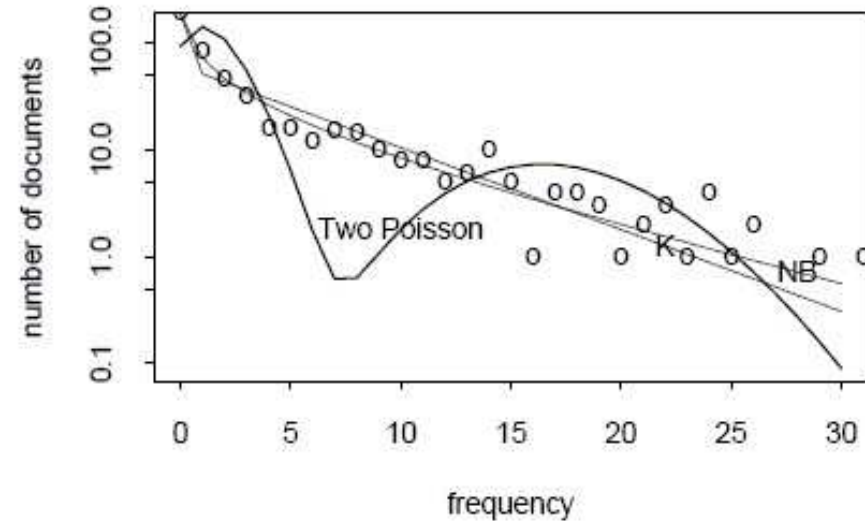
with $\Delta = \text{avgdoclength}$ and tuning parameters k_1, k_2, k_3, b , and non-linear influence of tf and consideration of doc length

Poisson Mixtures for Capturing tf Distribution

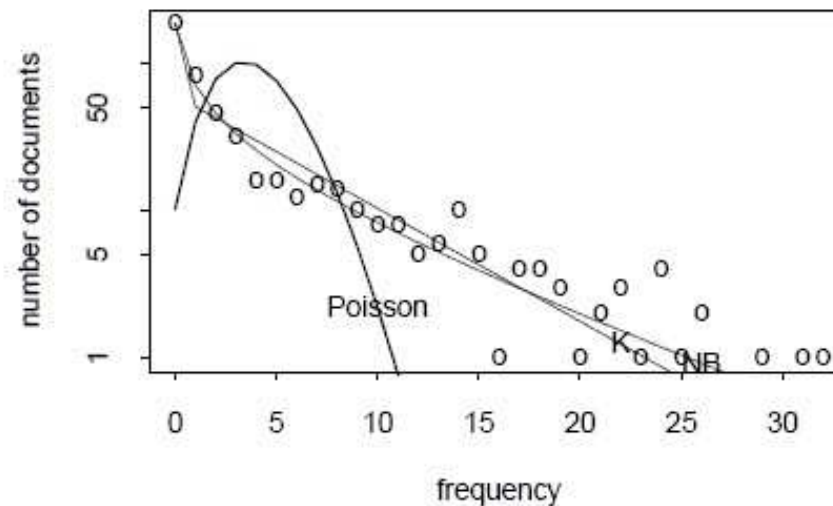
Poisson Doesn't Fit



Two Poissons Are Not Enough



Katz's K-mixture: Poisson Mixtures Fit Better



*distribution of
tf values
for term „said“*

Source:
Church/Gale 1995

Katz's K-Mixture

Katz's K-mixture: $f(k) = \int_0^{\infty} \Phi(\theta) \cdot \frac{e^{-\theta} \theta^k}{k!}$

e.g. with :

$$\Phi_K(\theta) = (1 - \alpha)\delta(\theta = 0) + \frac{\alpha}{\beta} e^{-\theta/\beta}$$

$$\rightarrow f(k) = (1 - \alpha)\delta(k = 0) + \frac{\alpha}{\beta + 1} \left(\frac{\beta}{\beta + 1} \right)^k$$

with $\delta(G)=1$ if G is true, 0 otherwise

Parameter estimation for given term:

$$\lambda = cf / N$$

observed mean tf

$$idf = \log_2(N / df)$$

$$\beta = \lambda 2^{idf} - 1 = (cf - df) / df$$

extra occurrences (tf>1)

$$\alpha = \lambda / \beta$$

4.1.4 Extensions of Probabilistic IR

Consider term correlations in documents (with binary X_i)

→ Problem of estimating m-dimensional prob. distribution

$$P[X_1=\dots \wedge X_2= \dots \wedge \dots \wedge X_m=\dots] =: f_{\mathbf{X}}(X_1, \dots, X_m)$$

One possible approach: **Tree Dependence Model:**

a) Consider only 2-dimensional probabilities (for term pairs)

$$f_{ij}(X_i, X_j) = P[X_i=\dots \wedge X_j=\dots] = \sum_{X_1} \dots \sum_{X_{i-1}} \sum_{X_{i+1}} \dots \sum_{X_{j-1}} \sum_{X_{j+1}} \dots \sum_{X_m} P[X_1 = \dots \wedge \dots \wedge X_m = \dots]$$

b) For each term pair

estimate the error between independence and the actual correlation

c) Construct a tree with terms as nodes and the

m-1 highest error (or correlation) values as weighted edges

Considering Two-dimensional Term Correlation

Variant 1:

Error of approximating f by g (**Kullback-Leibler divergence**)
with g assuming pairwise term independence:

$$\varepsilon(f, g) := \sum_{\vec{X} \in \{0,1\}^m} f(\vec{X}) \log \frac{f(\vec{X})}{g(\vec{X})} = \sum_{\vec{X} \in \{0,1\}^m} f(\vec{X}) \log \frac{f(\vec{X})}{\prod_{i=1}^m g_i(X_i)}$$

Variant 2:

Correlation coefficient for term pairs:

$$\rho(X_i, X_j) := \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}(X_i)} \sqrt{\text{Var}(X_j)}}$$

Variant 3:

level- α values or p-values
of **Chi-square independence test**

Example for Approximation Error ε (KL Strength)

$m=2$:

given are documents:

$$d1=(1,1), d2=(0,0), d3=(1,1), d4=(0,1)$$

estimation of 2-dimensional prob. distribution f :

$$f(1,1) = P[X1=1 \wedge X2=1] = 2/4$$

$$f(0,0) = 1/4, f(0,1) = 1/4, f(1,0) = 0$$

estimation of 1-dimensional marginal distributions $g1$ and $g2$:

$$g1(1) = P[X1=1] = 2/4, g1(0) = 2/4$$

$$g2(1) = P[X2=1] = 3/4, g2(0) = 1/4$$

estimation of 2-dim. distribution g with independent X_i :

$$g(1,1) = g1(1)*g2(1) = 3/8,$$

$$g(0,0) = 1/8, g(0,1) = 3/8, g(1,0) = 1/8$$

approximation error ε (KL divergence):

$$\varepsilon = 2/4 \log 4/3 + 1/4 \log 2 + 1/4 \log 2/3 + 0$$

Constructing the Term Dependence Tree

Given:

complete graph (V, E) with m nodes $X_i \in V$ and m^2 undirected edges $\in E$ with weights ε (or ρ)

Wanted:

spanning tree (V, E') with maximal sum of weights

Algorithm:

Sort the m^2 edges of E in descending order of weight

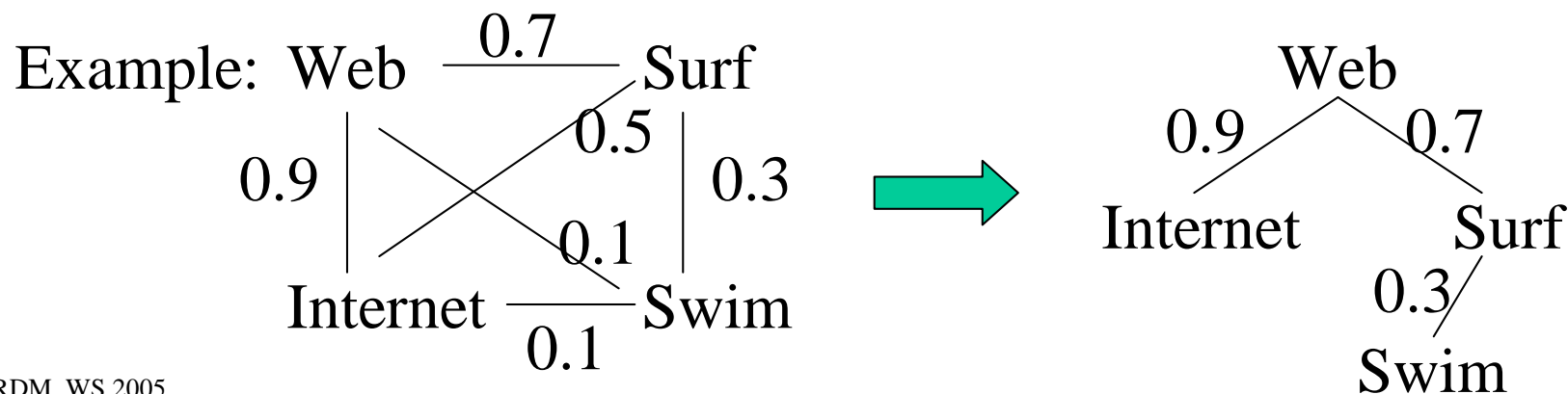
$E' := \emptyset$

Repeat until $|E'| = m-1$

$E' := E' \cup \{(i,j) \in E \mid (i,j) \text{ has max. weight in } E\}$

provided that E' remains acyclic;

$E := E - \{(i,j) \in E \mid (i,j) \text{ has max. weight in } E\}$



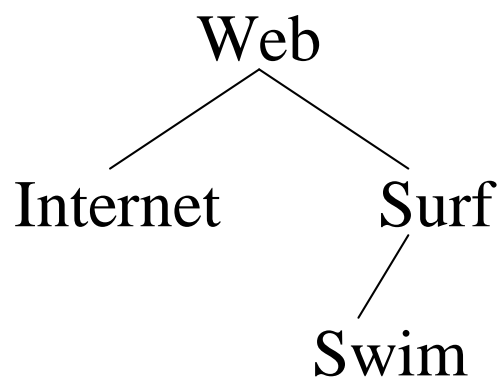
Estimation of Multidimensional Probabilities with Term Dependence Tree

Given is a term dependence tree $(V = \{X_1, \dots, X_m\}, E')$.

Let X_1 be the root, nodes are preorder-numbered, and assume that X_i and X_j are independent for $(i, j) \notin E'$. Then:

$$\begin{aligned}
 P[X_1 = .. \wedge .. \wedge X_m = ..] &= P[X_1 = ..] P[X_2 = .. \wedge X_m = .. | X_1 = ..] \\
 &= \prod_{i=1..m} P[X_i = .. | X_1 = .. \wedge X_{(i-1)} = ..] \\
 &= P[X_1] \cdot \prod_{(i,j) \in E'} P[X_j | X_i] \\
 &= P[X_1] \cdot \prod_{(i,j) \in E'} \frac{P[X_i, X_j]}{P[X_i]}
 \end{aligned}$$

Example:



$P[\text{Web}, \text{Internet}, \text{Surf}, \text{Swim}] =$

$$P[\text{Web}] \frac{P[\text{Web}, \text{Internet}]}{P[\text{Web}]} \frac{P[\text{Web}, \text{Surf}]}{P[\text{Web}]} \frac{P[\text{Surf}, \text{Swim}]}{P[\text{Surf}]}$$

Bayesian Networks

A **Bayesian network (BN)** is a directed, acyclic graph (V, E) with the following properties:

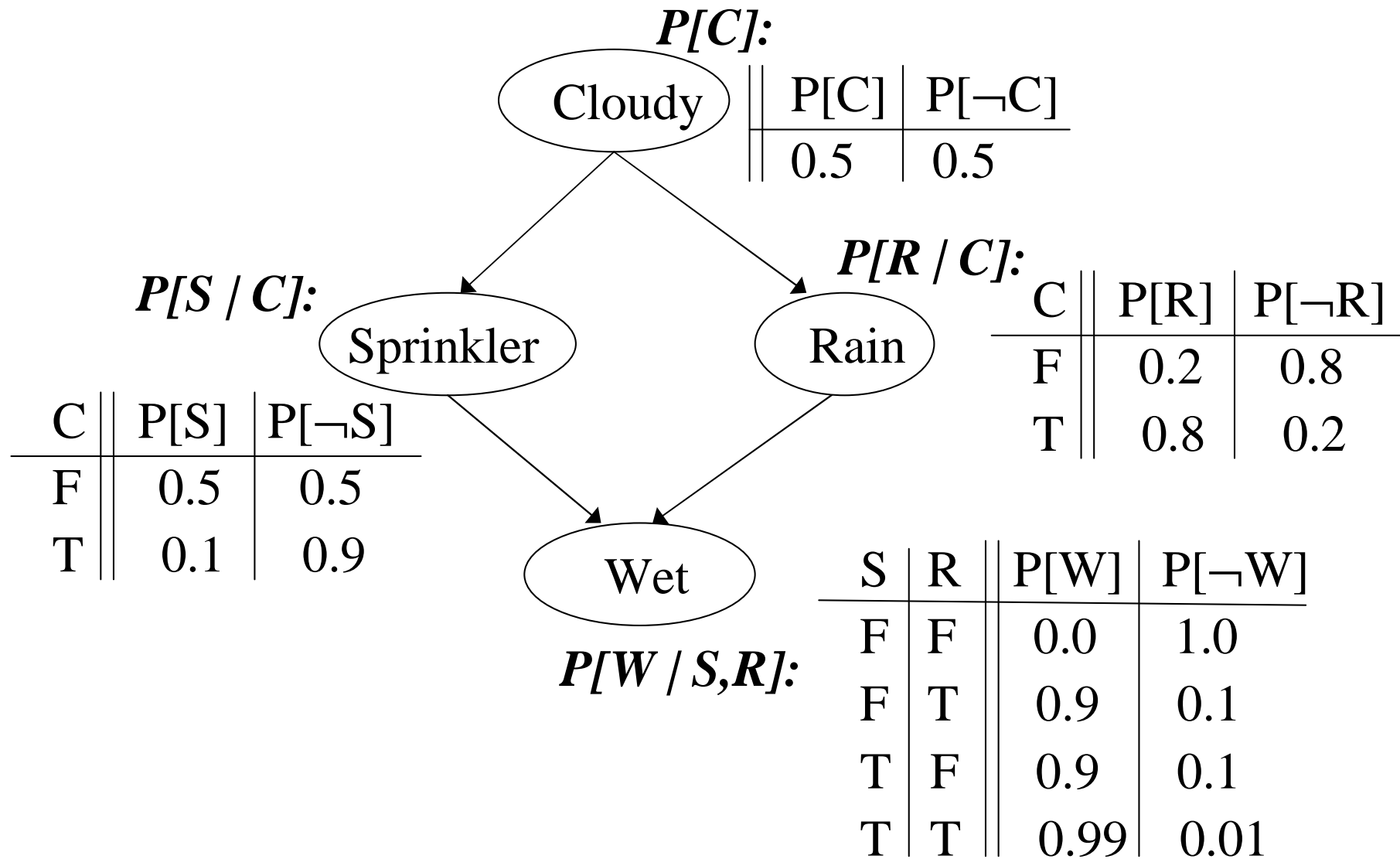
- Nodes $\in V$ representing random variables and
- Edges $\in E$ representing dependencies.
- For a root $R \in V$ the BN captures the prior probability $P[R = \dots]$.
- For a node $X \in V$ with parents $\text{parents}(X) = \{P_1, \dots, P_k\}$ the BN captures the conditional probability $P[X = \dots \mid P_1, \dots, P_k]$.
- Node X is conditionally independent of a non-parent node Y given its parents $\text{parents}(X) = \{P_1, \dots, P_k\}$:
 $P[X \mid P_1, \dots, P_k, Y] = P[X \mid P_1, \dots, P_k]$.

This implies: $P[X_1 \dots X_n] = P[X_1 \mid X_2 \dots X_n] P[X_2 \dots X_n]$

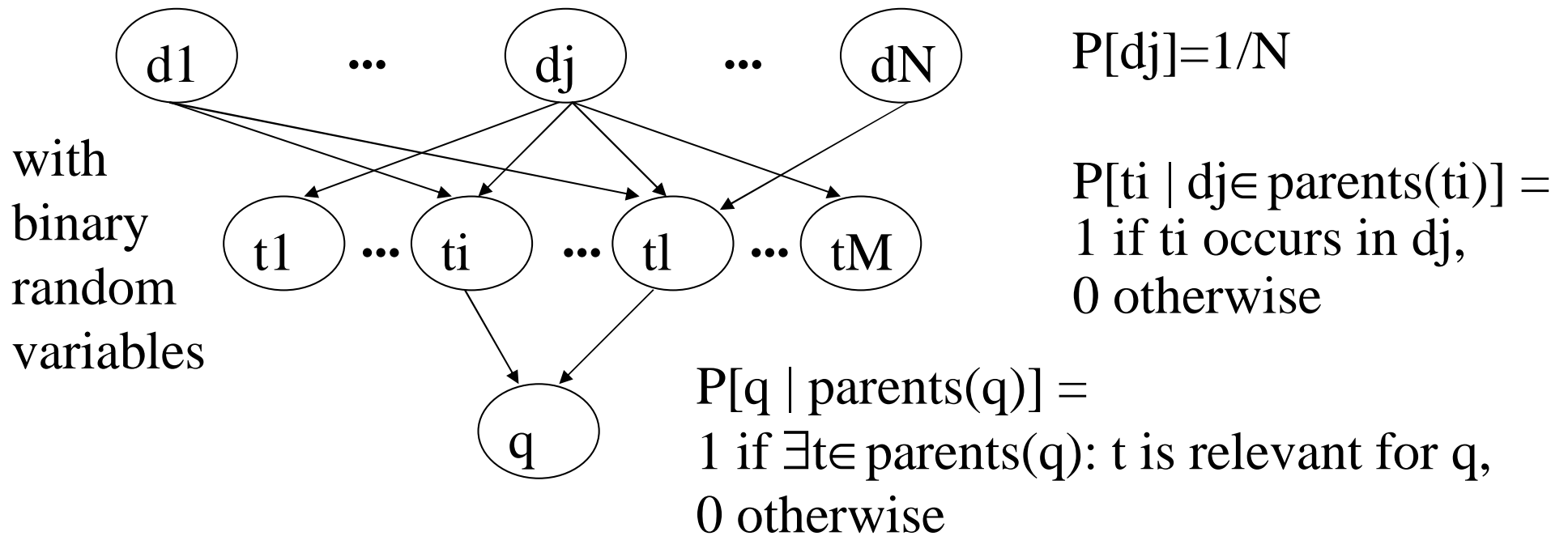
- by the chain rule:
$$= \prod_{i=1}^n P[X_i \mid X_{(i+1)} \dots X_n]$$
- by cond. independence:
$$= \prod_{i=1}^n P[X_i \mid \text{parents}(X_i), \text{other nodes}]$$

$$= \prod_{i=1}^n P[X_i \mid \text{parents}(X_i)]$$

Example of Bayesian Network (Belief Network)

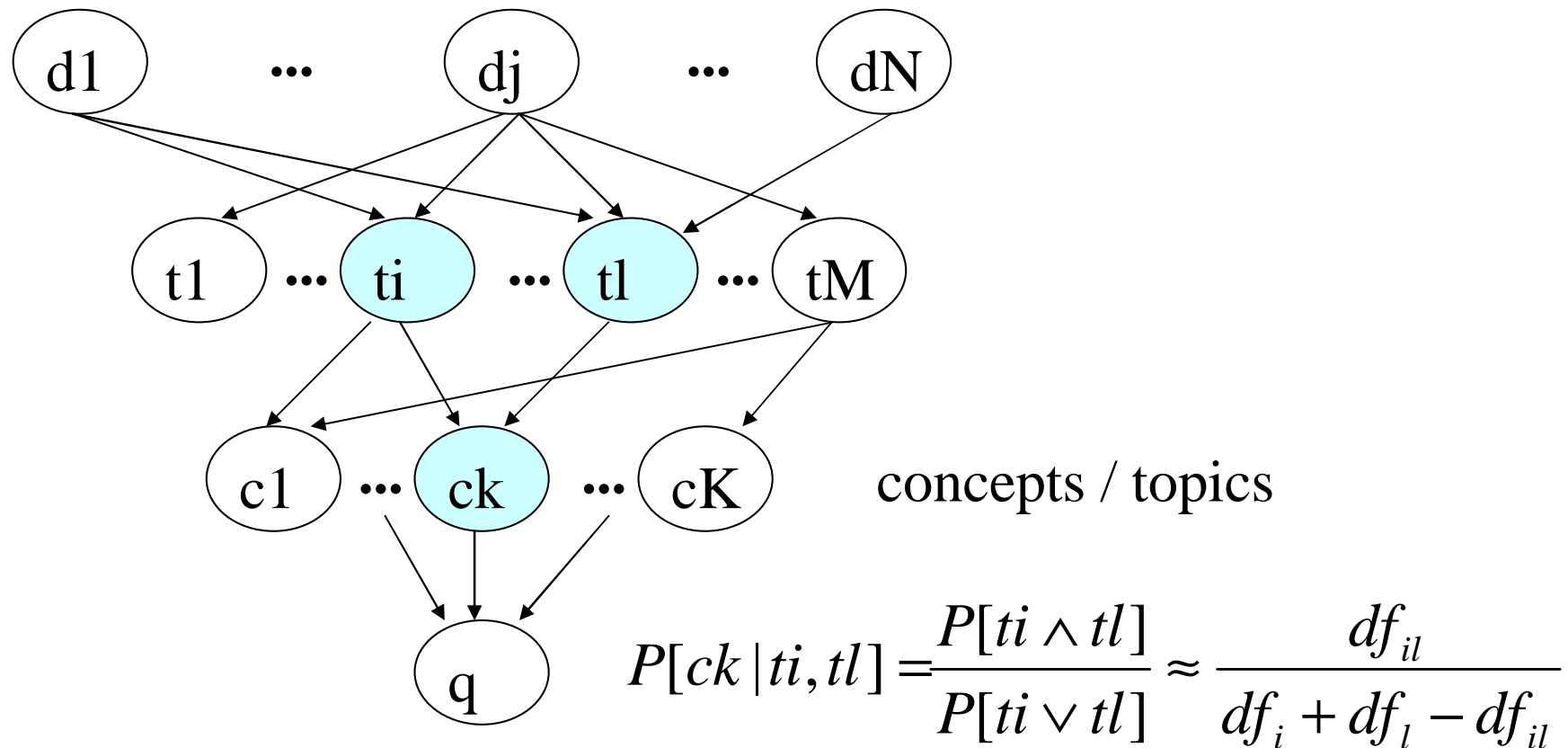


Bayesian Inference Networks for IR



$$\begin{aligned}
 P[q \wedge d_j] &= \sum_{(t_1 \dots t_M)} P[q \wedge d_j / t_1 \dots t_M] P[t_1 \dots t_M] \\
 &= \sum_{(t_1 \dots t_M)} P[q \wedge d_j \wedge t_1 \wedge \dots \wedge t_M] \\
 &= \sum_{(t_1 \dots t_M)} P[q / d_j \wedge t_1 \wedge \dots \wedge t_M] P[d_j \wedge t_1 \wedge \dots \wedge t_M] \\
 &= \sum_{(t_1 \dots t_M)} P[q / t_1 \wedge \dots \wedge t_M] P[t_1 \wedge \dots \wedge t_M / d_j] P[d_j]
 \end{aligned}$$

Advanced Bayesian Network for IR



Problems:

- parameter estimation (sampling / training)
- (non-) scalable representation
- (in-) efficient prediction
- fully convincing experiments

Additional Literature for Chapter 4

Probabilistic IR:

- Grossman/Frieder Sections 2.2 and 2.4
- S.E. Robertson, K. Sparck Jones: Relevance Weighting of Search Terms, JASIS 27(3), 1976
- S.E. Robertson, S. Walker: Some Simple Effective Approximations to the 2-Poisson Model for Probabilistic Weighted Retrieval, SIGIR 1994
- K.W. Church, W.A. Gale: Poisson Mixtures, Natural Language Engineering 1(2), 1995
- C.T. Yu, W. Meng: Principles of Database Query Processing for Advanced Applications, Morgan Kaufmann, 1997, Chapter 9
- D. Heckerman: A Tutorial on Learning with Bayesian Networks, Technical Report MSR-TR-95-06, Microsoft Research, 1995