

MODEL PROBABILITAS

2 November 2006

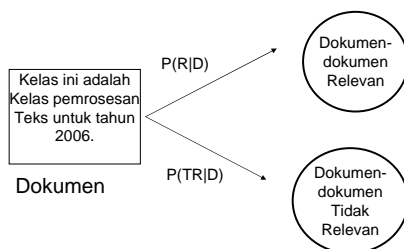
Bab 2: Grossman

Bab 7: Baeza-Yates & Ribeiro-Neto

Model Probabilitas

- Model probabilitas digunakan untuk memperoleh suatu set jawaban bagi suatu query.
- Bila ada suatu query, maka ada suatu set jawaban yang *ideal*
- Deskripsi dari suatu set jawaban ideal dimodelkan dalam bentuk probabilitas
 - Properti dari set jawaban ideal itu apa saja?

Model Probabilitas Dasar



Model Probabilitas Dasar

- Asumsikan relevansi suatu dokumen adalah independen dari dokumen-dokumen lain di dalam koleksi.
- Urutkan dokumen-dokumen sesuai nilai probabilitas relevansinya.

Prinsip Urutan Probabilitas

- Jika ada suatu query q dan suatu dokumen d_j , model probabilitas mencoba untuk memperkirakan probabilitas bahwa pemakai akan menemukan dokumen d_j relevan.
- Set jawaban yang ideal disebut sebagai R dan mestinya menunjukkan nilai probabilitas relevan yang maksimal. Dokumen-dokumen pada set R diramalkan sebagai relevan.
- Bagaimana cara menghitung probabilitas?

Fakultas Ilmu Komputer – Universitas Indonesia

MA-5

Model Probabilitas

- Tahapan penghitungan probabilitas dari suatu dokumen
 - Peroleh suatu set dokumen awal
 - Pemakai meneliti apakah dok ini adalah relevan (biasanya hanya 10-20 dok teratas)
 - Sistem IR menggunakan informasi ini untuk memperbaiki deskripsi dari set jawaban ideal
 - Dengan mengulangi proses ini, diharapkan deskripsi dari set jawaban ideal akan meningkat
- Pada awalnya perlu untuk menebak deskripsi dari suatu set jawaban ideal.

Fakultas Ilmu Komputer – Universitas Indonesia

MA-6

Beberapa Dasar Probabilitas

Bayes' Rule:

$$P(B|A) = P(A|B) * P(B) / P(A)$$

$$p(R | x) = \frac{p(x | R)p(R)}{p(x)}$$

$$p(NR | x) = \frac{p(x | NR)p(NR)}{p(x)}$$

Fakultas Ilmu Komputer – Universitas Indonesia

MA-7

Urutan

- Urutan dengan Probabilitas dihitung dengan:
 - $sim(q, dj) = P(dj \text{ relevant-to } q) / P(dj \text{ non-relevant-to } q)$
 - Ini adalah kemungkinan bahwa dokumen d_j adalah relevan
 - Dengan menggunakan kemungkinan akan memperkecil probabilitas dari suatu penilaian yang salah
- Definisi:
 - $P(R | vec(dj))$: probabilitas bahwa suatu dok relevan
 - $P(-R | vec(dj))$: probabilitas dok itu tidak relevan

Fakultas Ilmu Komputer – Universitas Indonesia

MA-8

Urutan

$$\begin{aligned} \bullet \text{ sim}(d_j, q) &= \frac{P(R \mid \text{vec}(d_j))}{P(\neg R \mid \text{vec}(d_j))} \\ &= \frac{[P(\text{vec}(d_j) \mid R) * P(R)]}{[P(\text{vec}(d_j) \mid \neg R) * P(\neg R)]} \\ &\sim \frac{P(\text{vec}(d_j) \mid R)}{P(\text{vec}(d_j) \mid \neg R)} \end{aligned}$$

- $P(\text{vec}(d_j) \mid R)$: probabilitas dari pemilihan dokumen secara acak d_j dari set dokumen yang relevan R

Urutan

$$\begin{aligned} \bullet \text{ sim}(d_j, q) &\sim \frac{P(\text{vec}(d_j) \mid R)}{P(\text{vec}(d_j) \mid \neg R)} \\ &\sim \frac{[\prod P(k_i \mid R)]}{[\prod P(k_i \mid \neg R)]} \end{aligned}$$

- $P(k_i \mid R)$: probabilitas kata indeks k_i ada di suatu dokumen yang dipilih secara acak dari set dokumen yang relevan R.

Urutan Awal

- Probabilitas $P(k_i \mid R)$ dan $P(k_i \mid \neg R)$?
- Perkiraan berdasarkan asumsi:
 - $P(k_i \mid R) = 0.5$
 - $P(k_i \mid \neg R) = \frac{n_i}{N}$
dimana n_i adalah jumlah dok yang berisi k_i
 - Gunakan perkiraan ini untuk memperoleh urutan awal
 - Lakukan perbaikan pada urutan awal

Memperbaiki Urutan Awal

- Asumsikan
 - V : set dari dok yang diperoleh pada awal proses
 - V_i : subset dari dok yang diperoleh yg berisi k_i
- Evaluasi kembali perkiraan:
 - $P(k_i \mid R) = \frac{V_i}{V}$
 - $P(k_i \mid \neg R) = \frac{n_i - V_i}{N - V}$
- Ulangi secara rekursif.

Contoh Model Probabilitas

- Nilai probabilitas kata pada dokumen dj

Kata	P(k R)	P(k -R)
k1	0.8	0.4
k2	0.6	0.1
k3	0.2	0.9
k4	0.9	0.6

$$P(R) = 0.1 \quad ; \quad P(-R) = 0.9$$

$$\text{Sim}(d_j, q) = \frac{0.8 * 0.6 * 0.2 * 0.9 * 0.1}{0.4 * 0.1 * 0.9 * 0.6 * 0.9} = \frac{0.00784}{0.01944} = 0.403$$