# Relevance Models

Berlin Chen 2005

References:
1. W. B. Croft and J. Lafferty (Editors). *Language Modeling for Information Retrieval*. Chapter 2, July 2003
2. V. Lavrenko and W.B. Croft, "Relevance-Based Language Models" ACM SIGIR 2001

# Introduction

- **Probabilistic Models** (e.g., Robertson and Sparck Jones, 1976)
  - Ranking docs by the odds of their being observed in the relevant class

$$SIM\ (Q,D) \approx \frac{P\left(D\middle|R_Q\right)}{P\left(D\middle|\overline{R}_Q\right)}$$

  - Explicitly model the class/concept of relevance

- Language Modeling Approaches
  - View documents themselves as models and queries as strings of text (observations) randomly sampled from these models

$$SIM\ (Q,D) \approx P\left(Q\middle|D\right) \approx \prod_{q_n \in Q} P\left(q_n\middle|D\right)$$ bag-of-words assumption

  - E.g., as the HMM, TMM, PLSA, mentioned previously

  - Quite several probability estimation techniques were proposed

# Relevance Models (RM)

- A convergence of ideas from classical probabilistic approaches and language modeling approaches
  - Explicit model the class $R$ of relevant documents

  - Describe unigram generative models for a set of text samples( in the relevant class" $P(w|R)$
    - Documents are randomly drawn from such a class

# RM: Ranking Approaches

- Probability Ratio $\left(\text{given that } D = d_1 \cdots d_n\right)$

$$\frac{P\left(D \mid R\right)}{P\left(D \mid N\right)} = \frac{P\left(d_1 \cdots d_n \mid R\right)}{P\left(d_1 \cdots d_n \mid N\right)} \approx \frac{\prod\limits_{i=1}^{n} P\left(d_i \mid R\right)}{\prod\limits_{i=1}^{n} P\left(d_i \mid N\right)} \approx \frac{\prod\limits_{i=1}^{n} P\left(d_i \mid R\right)}{\prod\limits_{i=1}^{n} P\left(d_i \mid C\right)}$$

- Higher score is better (more relevant)

---

- $R$ : the relevant class
- $N$ : the non-relevant class
- $C$ : the document collection or corpus

$$\left(|R| \text{is small} \Rightarrow |N| \approx |C|, \text{all measured in terms of word number}\right)$$

- $P\left(d_i \mid R\right)$ or $P\left(w \mid R\right)$ can be estimated from a set of training examples or from the query alone  (a focus of much active research)

---

(Inspired by the probability ranking principle)

# RM: Ranking Approaches (cont.)

- Cross-Entropy

$$P\left(R\|D\right) = -\sum_{w \in V} P\left(w|R\right)\log P\left(w|D\right)$$

<span style="color:blue">vocabulary</span>

- – Lower score ( $\geq 0$ ) is better (more relevant)

- – The negation of this ranking approach is equivalent to "query-likelihood" language-modeling approach (a special case of the cross-entropy approach itself ?)
  - Suppose $P\left(w|R\right)$ is the relative frequency of $w$ in the query $Q$

$$- \quad P\left(R\|D\right) \approx \sum_{w \in V} P\left(w|Q\right)\log P\left(w|D\right) = \sum_{w \in V} \frac{C\left(w,Q\right)}{|Q|}\log P\left(w|D\right)$$

$$= \frac{1}{|Q|}\sum_{w \in V}\log P\left(w|D\right)^{C\left(w,Q\right)} = \frac{1}{|Q|}\log \prod_{w \in V} P\left(w|D\right)^{C\left(w,Q\right)}$$

$$= \frac{1}{|Q|}\log P\left(Q|D\right)$$

# RM: Estimation from a Set of Examples

- Given that we have perfect knowledge of the entire relevant class $R$

$$P(w|R) = \sum_{D \in Collection} P(w, D|R)$$

$$= \sum_{D \in Collection} P(w|D, R) P(D|R)$$

conditional independence assumption

$$= \sum_{D \in Collection} P(w|D) P(D|R)$$

$$= \sum_{D \in Collection} P_{ML}(w|D) P(D|R)$$

- Where

no. of documents in $R$

$$P_{ML}(w|D) = \frac{C(w, D)}{|D|} \qquad P(D|R) = \begin{cases} 1/|R| & \text{if } D \in R \\ 0 & \text{otherwise} \end{cases}$$

# RM: Estimation from a Set of Examples (cont.)

- Probability smoothing by simple interpolation

$$P_{smooth}\left(w|D\right) = \lambda_D\, P_{ML}\left(w|D\right) + \left(1 - \lambda_D\right)P\left(w\right)$$

where

no. of documents in the collection

$$P\left(w\right) = \sum_{D \in C} P_{ML}\left(w|D\right) \cdot \frac{1}{|C|}$$

- Which is also equivalent to adjust document posterior probability

$$P_{smooth}\left(D|R\right) = \begin{cases} \dfrac{\left(1 - \lambda_D\right)}{|C|} + \dfrac{\lambda_D}{|R|} & \text{if } D \in R \\[2ex] \dfrac{\left(1 - \lambda_D\right)}{|C|} & \text{otherwise} \end{cases}$$