

O modelo probabilístico é baseado no princípio da ordenação probabilística (*Probability Ranking Principle*) onde dada uma consulta  $q$  e um documento  $d_j$  [perfeitamente] relevante a  $q$ , o modelo tenta estimar a probabilidade do usuário encontrar o documento  $d_j$ . O modelo assume que para uma consulta  $q$  há um conjunto de documentos  $R$  que contém exatamente os documentos relevantes e nenhum outro, sendo este um conjunto resposta ideal que maximiza a probabilidade do usuário encontrar um documento  $d_j$  relevante a  $q$ .

Seja  $\bar{R}_q$  o complemento de  $R$  de forma que  $\bar{R}_q$  contém todos os documentos não relevantes à consulta  $q$ . Seja  $P(R_q|d_j)$  a probabilidade do documento  $d_j$  ser relevante à consulta  $q$  e  $P(\bar{R}_q|d_j)$  a probabilidade de  $d_j$  não ser relevante à  $q$ . A similaridade entre um documento  $d_j$  e uma consulta  $q$  é definida por:

$$sim(d_j, q) = \frac{P(R_q|d_j)}{P(\bar{R}_q|d_j)} \quad (1)$$