

Language Models

- Consider the DIR model (query likelihood with Dirichlet smoothing slide 36 "Probabilistic IR")

$$P(q|\theta_d) = \frac{|q|!}{\prod_{t \in V} (|w_t^q|!)^{|q|}} \prod_{t \in V} p(t|\theta_d)^{w_t^q} \propto \prod_{t \in V} p(t|\theta_d)^{w_t^q}$$

- When using logs:
$$P(q|\theta_d) \propto \ln(\prod_{t \in V} p(t|\theta_d)^{w_t^q})$$
$$\propto \sum_{t \in V} \ln(p(t|\theta_d)^{w_t^q})$$
$$\propto \sum_{t \in V} w_t^q \cdot \ln(p(t|\theta_d))$$

Dirichlet Smoothing Language (DIR) Model

- Theoretical formula, with $\hat{\theta}_d$ the Dirichlet smoothed document language model ($\mu \in \mathbb{R}^+$)

$$P(q | \hat{\theta}_d) \propto \sum_{w \in V} c(w, q) \cdot \ln\left(\frac{c(w, d) + \mu \cdot P(w | C)}{|d| + \mu}\right)$$

- Good from a theoretical point of view, but not from an efficiency point of view as, as not directly implementable using inverted files
- Actual implementation using inverted files :

$$f(d, q) = \sum_{w \in d \cap q} c(w, q) \cdot \ln\left(1 + \frac{c(w, d)}{\mu \cdot P(w | C)}\right) + |q| \cdot \ln\left(\frac{\mu}{\mu + |d|}\right)$$

... Not obvious... explanations follow...

DIR model

- If $\hat{\theta}_d$ is the Dirichlet smoothed document language model, $\mu \in \mathbb{R}^+$

$$\begin{aligned}\ln P(q|\hat{\theta}_d) &=_{rank} \sum_{w \in V} c(w,q) \cdot \ln\left(\frac{c(w,d) + \mu \cdot P(w|C)}{|d| + \mu}\right) \\ &=_{rank} \sum_{w \in d \cap q} c(w,q) \cdot \ln\left(\frac{c(w,d) + \mu \cdot P(w|C)}{|d| + \mu}\right) + \sum_{w \in q, w \notin d} c(w,q) \cdot \ln\left(\frac{\mu \cdot P(w|C)}{|d| + \mu}\right) \\ &=_{rank} \sum_{w \in d \cap q} c(w,q) \cdot \ln\left(\frac{c(w,d) + \mu \cdot P(w|C)}{|d| + \mu}\right) + \sum_{w \in q, w \notin d} c(w,q) \cdot \ln\left(\frac{\mu \cdot P(w|C)}{|d| + \mu}\right) \\ &\quad + \sum_{w \in q} c(w,q) \cdot \ln\left(\frac{\mu \cdot P(w|C)}{|d| + \mu}\right) - \sum_{w \in q} c(w,q) \cdot \ln\left(\frac{\mu \cdot P(w|C)}{|d| + \mu}\right)\end{aligned}$$

DIR model

$$\begin{aligned}\ln P(q|\hat{\theta}_d) &=_{rank} \sum_{w \in d \cap q} c(w,q) \cdot \ln\left(\frac{c(w,d) + \mu \cdot P(w|C)}{|d| + \mu}\right) - \sum_{w \in d \cap q} c(w,q) \cdot \ln\left(\frac{\mu \cdot P(w|C)}{|d| + \mu}\right) \\ &\quad + \sum_{w \in q} c(w,q) \cdot \ln\left(\frac{\mu \cdot P(w|C)}{|d| + \mu}\right) \\ &=_{rank} \sum_{w \in d \cap q} c(w,q) \cdot \ln\left(1 + \frac{c(w,d)}{\mu \cdot P(w|C)}\right) + \sum_{w \in q} c(w,q) \cdot \ln\left(\frac{\mu}{|d| + \mu}\right) \\ &\quad + \sum_{w \in q} c(w,q) \cdot \log P(w|C) \\ &=_{rank} \sum_{w \in d \cap q} c(w,q) \cdot \ln\left(1 + \frac{c(w,d)}{\mu \cdot P(t|C)}\right) + |q| \cdot \ln\left(\frac{\mu}{|d| + \mu}\right)\end{aligned}$$