

1 计算entities概率的函数：getEntsPrblts()

标准式子：

$$p(y_d = l | w, z_{-d}, t_{-d}, y_{-d}, \Gamma) \propto \frac{\alpha_l + |\{y_{d'} = l, d \neq d'\}|}{\sum_{k=1}^K (\alpha_k + \alpha_{df} + |D| - 1)}$$

$$\prod_{i=1}^{N_d} \left(\frac{\gamma_1 + |\{t_{d'j} = 0, d \neq d'\}|}{\gamma_1 + \gamma_2 + \sum_{d \neq d'} N_{d'}} \frac{\beta_{bg} + |\{t_{d'j} = 0, w_{d'j} = w_{di}, d \neq d'\}|}{|W| \beta_{bg} + |\{t_{d'j} = 0, d \neq d'\}|} \right. \quad (1)$$

$$\left. + \frac{\gamma_2 + |\{t_{d'j} = 1, d \neq d'\}|}{\gamma_1 + \gamma_2 + \sum_{d \neq d'} N_{d'}} \frac{\beta_l + |\{z_{d'j} = l, w_{d'j} = w_{di}, d \neq d'\}|}{|W| \beta_l + |\{z_{d'j} = l, d \neq d'\}|} \right)$$

因为在同一个句子中许多变量都是相同的，在归一化是将被约去，所以简化了一部分分母：

$$p(y_d = l | w, z_{-d}, t_{-d}, y_{-d}, \Gamma) \propto (\alpha_l + |\{y_{d'} = l, d \neq d'\}|) \cdot$$

$$\prod_{i=1}^{N_d} \left((\gamma_1 + |\{t_{d'j} = 0, d \neq d'\}|) \cdot \frac{\beta_{bg} + |\{t_{d'j} = 0, w_{d'j} = w_{di}, d \neq d'\}|}{\beta_{bg} + \frac{|\{t_{d'j}=0, d \neq d'\}|}{|W|}} \right. \quad (2)$$

$$\left. + (\gamma_2 + |\{t_{d'j} = 1, d \neq d'\}|) \cdot \frac{\beta_l + |\{z_{d'j} = l, w_{d'j} = w_{di}, d \neq d'\}|}{\beta_l + \frac{|\{z_{d'j}=l, d \neq d'\}|}{|W|}} \right)$$

亦即：

$$p(y_d = l | w, z_{-d}, t_{-d}, y_{-d}, \Gamma) \propto (\alpha_l + |\{y_{d'} = l, d \neq d'\}|) \cdot$$

$$\prod_{i=1}^{N_d} \left(\frac{\gamma_1 + |\{t_{d'j} = 0, d \neq d'\}|}{\gamma_2 + |\{t_{d'j} = 1, d \neq d'\}|} \frac{\beta_{bg} + |\{t_{d'j} = 0, w_{d'j} = w_{di}, d \neq d'\}|}{\beta_{bg} + \frac{|\{t_{d'j}=0, d \neq d'\}|}{|W|}} \right. \quad (3)$$

$$\left. + \frac{\beta_l + |\{z_{d'j} = l, w_{d'j} = w_{di}, d \neq d'\}|}{\beta_l + \frac{|\{z_{d'j}=l, d \neq d'\}|}{|W|}} \right)$$

更显式的写为:

$$p(y_d = ent) \propto (\alpha_{ent} + Cnt_{ent(w.r.t \ pseudoDocs)}) \cdot \prod_{w \in pseudoDoc} \left(\frac{\gamma_{bg} + Cnt_{bg}}{\gamma_{df} + Cnt_{df}} \cdot \frac{\beta_{bg} + Cnt(bg, w)}{|W| * \beta_{bg} + Cnt_{bg}} + \frac{\beta_{ent} + Cnt(ent, w)}{|W| * \beta_{ent} + Cnt_{ent(w.r.t \ words)}} \right) \quad (4)$$

抽取一部分共同的式子:

$$comPart = (\gamma_{bg} + Cnt_{bg}) / (\gamma_{df} + Cnt_{df}) / (\beta_{bg} + Cnt_{bg} / |W|) \quad (5)$$

得:

$$p(y_d = ent) \propto (\alpha_{ent} + Cnt_{ent(w.r.t \ pseudoDocs)}) \cdot \prod_{w \in pseudoDoc} \left(comPart \cdot (\beta_{bg} + Cnt(bg, w)) + \frac{\beta_{ent} + Cnt(ent, w)}{\beta_{ent} + Cnt_{ent(w.r.t \ words)} / |W|} \right) \quad (6)$$

对此处的 $|w|$ 的具体所指存疑:

$$comPart = (\gamma_{bg} + Cnt_{bg}) / (\gamma_{df} + Cnt_{df}) / \left(\beta_{bg} + \frac{Cnt_{bg}}{|W|_{(:contextNum?currentContextNum)}} \right) \quad (7)$$

注: 这里的words都视为context words。

$$p(y_d = ent) \propto (\alpha_{ent} + Cnt_{ent}(w.r.t \ pseudoDocs)) \cdot \prod_{w \in pseudoDoc} \left(comPart \cdot (\beta_{bg} + Cnt(bg, w)) + \frac{\beta_{ent} + Cnt(ent, w)}{\beta_{ent} + \frac{Cnt_{ent}(w.r.t \ words)}{|W|(:contextNum?entity.wordid_{cnt.size()})}} \right) \quad (8)$$

ANS:context_words.

另外对于mention词，t必为1，所以对于此项word只需求：

$$\frac{\beta_{ent} + Cnt(ent, w)}{\beta_{ent} + Cnt_{ent}(w.r.t \ words)/|W|} \quad (9)$$

2 选择进一步增量学习数据的函数：isToLearn()

selectMode == 3时，使用置信度评定($> \eta$):

$$\frac{\max_i p(y_d = i)}{\sum_i p(y_d = i)} > \eta \quad (10)$$

$$\frac{\max_i e^{p(y_d=i)}}{\sum_i e^{p(y_d=i)}} > \eta \quad (11)$$

3 对unlabeled used set进行采样，sampleOnce(indexes):

3.1 问题是是否要对labeled data重采Indicators。

ANS:不用。

4 计算Indicator的函数：getIndicatorPrblt()

标准式子：

$$\frac{p(t_{di} = 0|w, z_{\neg d}, t_{\neg d}, y, \Gamma)}{p(t_{di} = 1|w, z_{\neg d}, t_{\neg d}, y, \Gamma)} = \frac{\gamma_1 + |\{t_{d'j} = 0, d' \neq d\}|}{\gamma_2 + |\{t_{d'j} = 1, d' \neq d\}|} \cdot \frac{\beta_{bg} + |\{t_{d'j} = 0, w_{d'j} = w_{di}, d' \neq d\}|}{\beta_{df} + |\{z_{d'j} = y_d, w_{d'j} = w_{di}, d' \neq d\}|} \cdot \frac{|W|\beta_{df} + |\{z_{d'j} = y_d, d' \neq d\}|}{|W|\beta_{bg} + |\{t_{d'j} = 0, d' \neq d\}|} \quad (12)$$

更显式的写为：

$$\frac{p(t_{di} = 0)}{p(t_{di} = 1)} = \frac{\text{gamma_bg} + \text{cnt_bg}}{\text{gamma_df} + \text{cnt_df}} \cdot \frac{\text{beta_bg} + \text{cnt}(\text{bg}, w_{di})}{\text{beta_df} + \text{cnt}(\text{ent}, w_{di})} \cdot \frac{|W|\text{beta_df} + \text{cnt_ent}}{|W|\text{beta_bg} + \text{cnt_bg}} \quad (13)$$