

6(a). When  $tf_{ij} > 0$ , what are the maximum and minimum values of the corresponding  $tf'_{ij}$  values based on the above transformation? Please specify what cases the max and min value achieves.

**Given:**

- $tf'_{ij} = tf_{ij} * (1 + \sum_0^j \frac{p_{ij} \log p_{ij}}{\log m})$  where  $p_{ij} = \frac{tf_{ij}}{gf_i}$
- $0 < tf_{ij} \leq k$
- $gf_i = \sum_1^m tf_{ij}$   
-above represents the number of times ith term appears in all documents
- $1 < m$   
-above represents max number of documents

**Proof:**

**Min:**

-Occurs when  $tf_{ij} = 1$  (aka only 1 of ith term found in each document)

Step 1: Assume  $tf_{ij} = 1$  and solve for  $p_{ij}$

$$p_{ij} = \frac{tf_{ij}}{\sum_1^m tf_{ij}} = \frac{1}{\sum_1^m 1} = \frac{1}{m}$$

Step 2: plug in the values for  $p_{ij}$  and  $tf_{ij}$  into the  $tf'_{ij}$  formula

$$\begin{aligned} tf'_{ij} &= tf_{ij} * (1 + \sum_0^j \frac{p_{ij} \log p_{ij}}{\log m}) \\ &= 1 * (1 + \sum_0^j \frac{\frac{1}{m} \log \frac{1}{m}}{\log m}) \\ &= 1 + j \left[ \frac{\frac{1}{m} \log \frac{1}{m}}{\log m} \right] \\ &= 1 + \frac{-j}{m} \\ &= \frac{m-j}{m} \end{aligned}$$

$$\mathbf{Ans:} = \frac{m-j}{m}$$

**Max:**

-Occurs when  $tf_{ij} = k$  (aka only k(max frequency) of ith term is found in each document)

Step 1: Assume  $tf_{ij} = k$  and solve for  $p_{ij}$

$$p_{ij} = \frac{tf_{ij}}{\sum_1^m tf_{ij}} = \frac{k}{\sum_1^m k} = \frac{k}{km} = \frac{1}{m}$$

Step 2: plug in the values for  $p_{ij}$  and  $tf_{ij}$  into the  $tf'_{ij}$  formula

$$\begin{aligned}
 tf'_{ij} &= tf_{ij} * (1 + \sum_0^j \frac{p_{ij} \log p_{ij}}{\log m}) \\
 &= k * (1 + \sum_0^j \frac{\frac{1}{m} \log \frac{1}{m}}{\log m}) \\
 &= k * (1 + \frac{-j}{m}) \\
 &= \frac{k(m-j)}{m}
 \end{aligned}$$

**Ans:**  $= \frac{k(m-j)}{m}$

b.) Briefly explain the purpose for this transformation in the context of natural language processing.

**Ans:** The purpose of this transformation is to get the average occurrence of the  $i$ th term found from the first to the  $j$ th document.