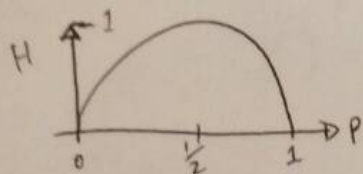


## HW # 2:

1) a)  $P(W = w_i) = \frac{1}{N}$

$$H(W) = - \sum_{w \in V} P(W) \log P(W)$$



$$H(W) = - \sum_{w \in V} \frac{1}{N} \log \frac{1}{N} = \sum_{w \in V} \frac{1}{N} \log N$$

$$= N \left( \frac{1}{N} \log N \right) = \log N$$

If the number of unique words  $N=1$  (i.e. all words are identical)  
then  $H(W) = \log 1 = 0$

$\therefore$  minimum  $H(W) = 0$  ; Maximum  $H(W) = \log N$

b) Sample Minimum  $H(W)$  article =  $\{w_1, w_1, w_1, w_1, w_1\}$   
(e.g. maximally homogeneous)

Sample maximum  $H(W)$  article =  $\{w_2, w_1, w_3, w_4, w_5, w_6, w_3, w_5\}$   
(e.g. maximally heterogeneous)

c) Two article which has  $H(W) = 0$  probably means that the documents themselves contain <sup>only</sup> one unique word each.

For example,  $D_1 = \{w_1, w_1, w_1, \dots, w_1\}$  and  $D_2 = \{w_2, w_2, \dots, w_2\}$

in that case, when combining the documents, the most distinct set that you would get is  $\{w_1, w_2\}$  so the maximum entropy for  $A_3$  is  $\log 2 = 0.69$