Q1: Which of these feature weighting techniques is best to give lower weights to words that occurs very frequency in all documents? (Give one line reason) [2 points]

I.      Binary occurrence
II.     Term Frequency
III.    Term count
IV.     TF-IDF

Q2: What is the underlying assumption in Vector space model? (Answer in 1 to 2 lines) [2 points]

Q3: Consider the following contingency table and retrieval status value (RSV) formula. * [2 points]

| | documents | relevant | nonrelevant | total |
|---|---|---|---|---|
| term present | $x_t = 1$ | $s$ | $df_t - s$ | $df_t$ |
| term absent | $x_t = 0$ | $S - s$ | $(N - df_t) - (S - s)$ | $N - df_t$ |
| | total | $S$ | $N - S$ | $N$ |

$$RSV_d = \sum_{t \in q} \left[ c_t \times \frac{(k_1 + 1)tf_{td}}{k_1((1-b) + b(L_d/L_{ave})) + tf_{td}} \times \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}} \right]$$

Where $c_t$ is

$$c_t = K(N, df_t, S, s) = \log \frac{s/(S-s)}{(df_t - s)/((N - df_t) - (S - s))}.$$

Change the above formula if we assume that all the document containing term t are relevant.

*$df_t$= document frequency of term t i.e. the number of document in which t occurs$tf_{td}$: Term Frequency of t in document d
$L_d$ : length of document d          $L_{ave}$= Average length of douments$tf_{tq}$: Term Frequency of term t in query q