# 💡 Precision and Recall:

$$P = \frac{m}{w_c}$$

$$R = \frac{m}{w_R}$$

→ $m$ = # of n-grams in candidates also found in reference.

→ $w_c$ = # of n-grams in candidate.

→ $w_r$ = # of n-grams in reference.

### EX:

R: I work on machine learning.

$C_A$: I work

$C_B$: He works on machine learning.

$$P_A = \frac{2}{2} = 100\%$$

$$P_B = \frac{3}{5} = 60\%$$

$$R_A = \frac{2}{5} = 40\%$$

$$R_B = \frac{3}{5} = 60\%$$

→ But the precision and recall are not always proper measure.

# 💡 BLEU:

→ Bilingual Evaluation Understudy.

→ $$BLEU = BP * \exp\left(\sum P_n\right)$$

→ BP = Brevity penalty

 → this adjusts the score when candidate is shorter than reference.

 → $BP = \begin{cases} 1 & \text{if } c > \gamma \\ \exp(1 - \frac{c}{\gamma}) & \text{otherwise} \end{cases}$

 → $P_n$ = n-gram precision

🟢 → easy to calculate and interpret

🔴 → heavily relies on n-grams which may not capture overall meaning or fluency.

🔴 → always penalizes when transalation (candidate) is longer than reference.

# 💡 ROUGE:

→ Recall oriented understudy for gisting evaluation.

→ mostly used for text summarization

→ based on Recall.

→ $ROUGE = \sum (\text{Recall of n-grams})$

→ there are multiple measure-

    1. ROUGE-N
    2. ROUGE-L
    3. ROUGE-S

## 2. ROUGE-L:

→ based on longest common subsequence.

→ $P = \dfrac{LCS(A,B)}{m}$     $R = \dfrac{LCS(A,B)}{n}$

    $m$ = candidate length
    $n$ = reference length
    $A$ = candidate
    $B$ = reference

then weighted harmonic mean ($F_1$):

$$F_1 = \frac{(1+b^2)RP}{R^2 + bP}$$

- ROUGE-W: weighted LCS
- ROUGE-S: skipgram allowed.

# 💡 METEOR:

→ having better correlation with human judgement.

→ $P = \dfrac{m}{W_c}$    $R = \dfrac{m}{W_R}$

$$F_{mean} = \dfrac{10\, PR}{R + 9P}$$

## chunk penalty:

$$\boxed{p = \gamma \cdot \left(\dfrac{c}{U_m}\right)^{\beta}}$$

- $0 < \gamma \leq 1$
- 

R: the cat sat on the mat

C: on the mat sat the cat

chunk: consecutive set of words in the candidate matching with the consecutive set of words in reference.

$c =$ # of chunks in the candidate
$U_m =$ length of the candidate.

here,   $U_m = 6$        → the cat
        $c = 3$  ⟷  the mat
        $\gamma = 0.5$      → on the
        $\beta = 3$

- METEOR $= F_{mean}(1-p)$