



Question Answering

Evaluating Summaries:
ROUGE



ROUGE (Recall Oriented Understudy for Gisting Evaluation)

Lin and Hovy 2003

- Intrinsic metric for automatically evaluating summaries
 - Based on BLEU (a metric used for machine translation)
 - Not as good as human evaluation (“Did this answer the user’s question?”)
 - But much more convenient
- Given a document D, and an automatic summary X:
 1. Have N humans produce a set of reference summaries of D
 2. Run system, giving automatic summary X
 3. What percentage of the bigrams from the reference summaries appear in X?

$$ROUGE - 2 = \frac{\sum_{s \in \{\text{RefSummaries}\}} \sum_{\text{bigrams } i \in S} \min(\text{count}(i, X), \text{count}(i, S))}{\sum_{s \in \{\text{RefSummaries}\}} \sum_{\text{bigrams } i \in S} \text{count}(i, S)}$$



A ROUGE example:

Q: “What is water spinach?”

Human 1: Water spinach is a green leafy vegetable grown in the tropics.

Human 2: Water spinach is a semi-aquatic tropical plant grown as a vegetable.

Human 3: Water spinach is a commonly eaten leaf vegetable of Asia.

- System answer: Water spinach is a leaf vegetable commonly eaten in tropical areas of Asia.

- ROUGE-2 =
$$\frac{3 + 3 + 6}{10 + 9 + 9} = 12/28 = .43$$

[illegible][illegible]