

* PART OF SPEECH TAGGING \Rightarrow MORPHOLOGICAL DISAMBIGUATION

- OPEN & CLOSED CLASSES:

- OPEN CLASSES: NEW ONES CAN BE CREATED ALL THE TIME EX) NOUN, VERB, ADJ, ADV
- CLOSED CLASSES HAVE A SMALL FIXED MEMBERSHIP EX) PREP, AUX, PRONOUN, FUNCTION WORDS.

- FOR POS TAGGING, A STANDARD SET OF TAGS MUST BE CHOSEN \rightarrow PENN TREEBANK TAGSET (45 TAGS) COMMONLY USED.

\Rightarrow POS TAGGING PROBLEM IS DETERMINING THE POS TAG FOR A PARTICULAR INSTANCE OF A WORD.

- RULE-BASED TAGGING: ASSIGN ALL POSSIBLE TAGS TO WORDS FROM THE DICTIONARY

- \rightarrow WRITE RULES BY HAND TO SELECTIVELY REMOVE TAGS
- \rightarrow LEAVES THE CORRECT TAG FOR EACH WORD.

• ENGLISH / ENGLISH TAGGING: RUN WORDS THROUGH FST MORPHOLOGICAL ANALYZER TO GET ALL POS \rightarrow APPLY NEGATIVE CONSTRAINTS.

- STOCHASTIC TAGGING: PROBABILISTIC SEQUENCE MODELS

• HIDDEN MARKOV MODEL TAGGING (HMM): SPECIAL CASE OF ^{BAYESIAN} ~~BAYES~~ INFERENCE

- USAGE: SPEECH RECOGNITION \rightarrow OBSERVED: ACOUSTIC SIGNAL, HIDDEN: WORDS
- HANDWRITING RECOGNITION \rightarrow OBSERVED: IMAGE, HIDDEN: WORDS
- POS TAGGING \rightarrow OBSERVED: WORDS, HIDDEN: POS
- MT \rightarrow OBSERVED: FOREIGN WORDS, HIDDEN: WORDS IN TARGET LANG

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(t_1^n / w_1^n) \approx \underset{t_1^n}{\operatorname{argmax}} \prod_{i=1}^n \underbrace{P(w_i | t_i)}_{\text{LIKELIHOOD}} \underbrace{P(t_i | t_{i-1})}_{\text{PRIOR}}$$

- TAG TRANSITION PROBABILITIES $p(t_i | t_{i-1}) = \frac{c(t_{i-1}, t_i)}{c(t_{i-1})}$

- WORD LIKELIHOOD PROBABILITIES $p(w_i | t_i) = \frac{c(t_i, w_i)}{c(t_i)}$

\Rightarrow A MARKOV CHAIN IS A SPECIAL CASE OF A WFST, IN WHICH THE INPUT SEQUENCE UNIQUELY DETERMINES WHICH STATES THE AUTOMATON WILL GO THROUGH.

\rightarrow HIDDEN MARKOV MODEL IS AN EXTENSION OF A MARKOV CHAIN IN WHICH THE INPUT SYMBOLS ARE NOT THE SAME AS THE STATES.

- VITERBI ALGORITHM: CREATE AN ARRAY WITH COLUMNS CORRESPONDING TO INPUTS AND ROWS CORRESPONDING TO POSSIBLE STATES \rightarrow SWEEP THROUGH ~~THE~~ USING TRANSITION PROBS & OBSERVATION PROBS \rightarrow STORE ONLY MAX PROB PATH TO EACH CELL.

- EVALUATION OF POS TAGGING:

- COMPARING WITH MANUALLY CODED "GOLD STANDARD"

* SYNTAX → THE WAY WORDS ARE ARRANGED TOGETHER.

- CONSTITUENCY: GROUPS OF WORDS MAY BEHAVE AS A ~~SINGLE~~ SINGLE UNIT OR PHRASE
⇒ A CONSTITUENT
- CONSTITUENTS CAN BE SHOWN TO BEHAVE IN SIMILAR WAYS WITH RESPECT TO THEIR INTERNAL STRUCTURE & OTHER UNITS IN THE LANGUAGE.
- CONTEXT-FREE GRAMMARS (CFG) = MATHEMATICAL SYSTEM FOR MODELLING CONSTITUENT STRUCTURE.

• TERMINAL SYMBOLS: CORRESPOND TO WORDS IN THE ~~LANGUAGE~~ LANGUAGE

• NON-TERMINAL SYMBOLS: EXPRESS CLUSTERS OR GENERALIZATIONS

• RULES: EQUATIONS THAT CONSIST OF A SINGLE NON-TERMINAL ON THE LEFT & ANY NUMBER OF TERMINALS & NON-TERMINALS ON THE RIGHT

ex) NP → Det Nominal

NP → Proper Noun

Nominal → Noun | Nominal Noun

} RULES FOR NOUN PHRASES.

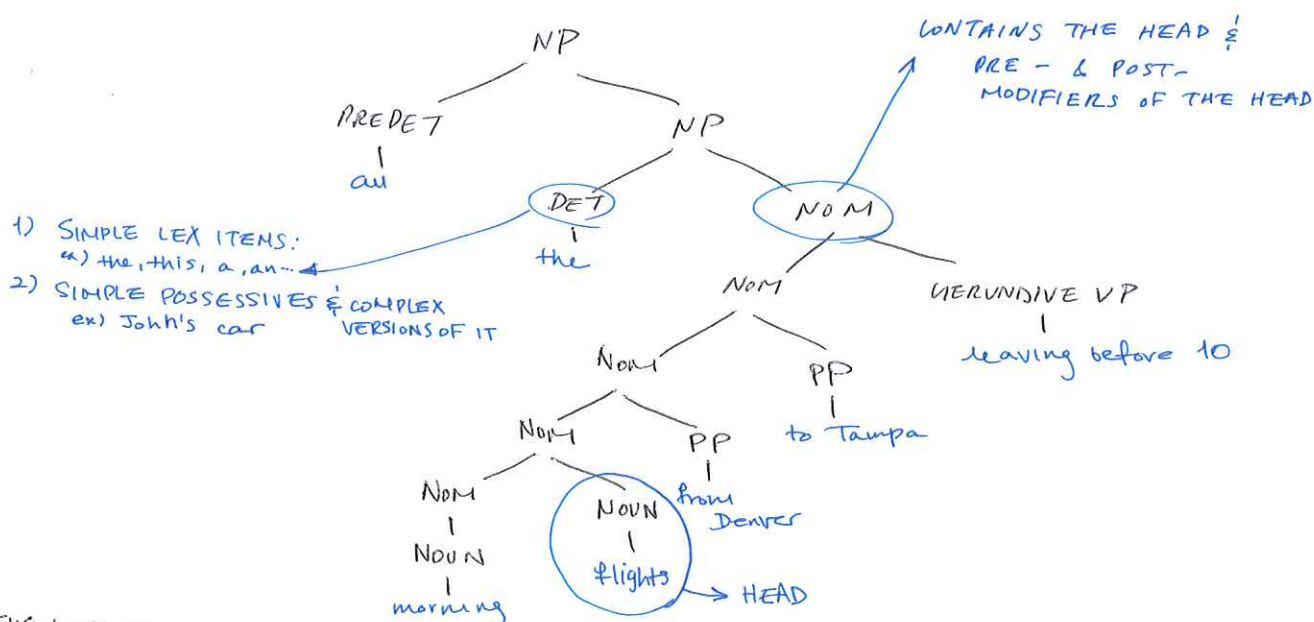
→ EXPLICIT DISJUNCTION - TWO KINDS OF NOMINALS

→ RECURSIVE DEFINITION - SAME NON-TERMINAL ON THE RIGHT & LEFT-SIDE OF THE RULE.

• DERIVATION: A SEQUENCE OF RULES APPLIED TO A STRING THAT ACCOUNTS FOR THAT STRING → COVERS ALL, BUT ONLY, THE ELEMENTS IN THE STRING

- SENTENCE TYPES: DECLARATIVES, IMPERATIVES, YES-NO QUESTIONS, WH-QUESTIONS.

- NOUN PHRASES:



→ THIS DERIVATION OVERGENERATES, AS IT DOES NOT CARE ABOUT AGREEMENTS

ex) DOES NOT DISTINGUISH B/W THIS FLIGHT & THESE FLIGHT

(ARGUMENTS)

- VERB PHRASES: CONSIST OF A HEAD VERB ALONG WITH 0 OR MORE FOLLOWING CONSTITUENTS

• SUBCATEGORIZE THE VERBS IN A LANGUAGE ACCORDING TO THE SETS OF VP RULES THAT THEY PARTICIPATE IN.

- CFGs PROVIDE INFORMATION ABOUT THE BASIC SYNTACTIC STRUCTURE BUT OVERGENERATE.

⇒ THERE ARE WAYS TO DEAL WITH THE PROBLEM, BUT NOT ELEGANT

⇒ THERE ARE SIMPLER & ELEGANT SOLUTIONS OUTSIDE OF THE CFG FRAMEWORK

* **TREEBANKS**: CORPORA IN WHICH EACH SENTENCE HAS BEEN PAIRED WITH A PARSE TREE

- IMPLICITLY DEFINE A GRAMMAR FOR THE LANGUAGE COVERED IN THE TREEBANK.
- THE DERIVED GRAMMAR IS VERY FLAT AS THEY TEND TO AVOID RECURSION

- **HEAD FINDING**: USE A SIMPLE SET OF TREE TRAVERSAL RULES SPECIFIC TO EACH NON-TERMINAL IN THE GRAMMAR

- **PARSING**: GIVEN A STRING OF TERMINALS & A CFG, DETERMINE IF THE STRING CAN BE GENERATED BY THE CFG → RETURN THE PARSE TREE(S) FOR THE STRING

→ **TOP-DOWN PARSING**: START SEARCHING SPACE OF DERIVATIONS FOR THE START SYMBOL

→ **BOTTOM-UP PARSING**: START SEARCH SPACE OF REVERSE DERIVATIONS FROM THE TERMINAL SYMBOLS IN THE STRING

→ NEVER EXPLORES OPTIONS THAT WILL NOT LEAD TO A FULL PARSE, BUT CAN EXPLORE ~~MANY~~ MANY OPTIONS THAT NEVER CONNECT TO THE ACTUAL SENTENCE

→ NEVER EXPLORES OPTIONS THAT DO NOT CONNECT TO THE ACTUAL SENTENCE, BUT CAN EXPLORE OPTIONS THAT CAN NEVER LEAD TO A FULL PARSE.

→ DYNAMIC ALGORITHMS BASED ON BOTH APPROACHES ACHIEVE $O(n^3)$ RECOGNITION TIME, WHERE n IS THE LENGTH OF THE INPUT STRING

→ CACHING IS CRITICAL TO OBTAINING A POLYNOMIAL TIME PARSING ALGORITHM FOR A CFG

- **CKY ALGORITHM**: PRODUCES ALL POSSIBLE PARSE TREES

• GRAMMAR MUST BE CONVERTED TO CHOMSKY NORMAL FORM (CNF) → LEXICAL RULES: EXACTLY 2 NON-TERMINAL SYMBOLS ON RHS OR 1 TERMINAL SYMBOL → MUST HAVE EITHER

• PARSE BOTTOM-UP STORING PHRASES FORMED FROM ALL SUBSTRINGS IN A TRIANGULAR TABLE (CHART)

• **COMPLEXITY**:

$$(n(n+1)/2) = O(n^2) \text{ CELLS} \times O(n) \text{ POSSIBLE SPLIT POINTS} = O(n^3)$$

→ SYNTACTIC AMBIGUITY CAN BE RESOLVED BY USING PROBABILISTIC CFG (PCFG) TO RUN CKY ALGORITHM TO DETERMINE THE MOST LIKELY PARSE TREE.

* **DEPENDENCY PARSING**

- **DEPENDENCY SYNTAX**: ASSUMES THAT SYNTACTIC STRUCTURE CONSISTS OF LEXICAL ITEMS LINKED BY BINARY ASYMMETRIC RELATIONS CALLED DEPENDENCIES

→ HEAD RULES CAN BE USED TO EXTRACT A DEPENDENCY PARSE FROM A CFG PARSE.

- **SOURCES OF INFORMATION FOR DEPENDENCY PARSING**:

• BILEXICAL AFFINITIES

• DEPENDENCY DISTANCE

• INTERVENING MATERIAL

• **VALENCY OF HEADS**: HOW MANY DEPENDENTS ON WHICH SIDE ARE USUAL FOR HEADS?

- **MatParser**: A SIMPLE FORM OF GREEDY DISCRIMINATIVE DEPENDENCY PARSER → DOES A SEQUENCE OF BOTTOM-UP ACTIONS

• MAINTAINS 4 STACKS:

- σ : STARTS WITH ROOT → WORD LEFT BETWEEN ROOT & . IS THE ROOT

- β : STARTS WITH THE INPUT SENTENCE → TERMINATES WHEN THIS IS EMPTY.

- A : A SET OF DEPENDENCY ARCS A → STARTS EMPTY

- SET OF ACTIONS

⇒ PROVIDES VERY FAST LINEAR TIME PARSING

- Projectivity: dependencies from a CFG tree using heads, must be projective.
 - there must not be any crossing dependency arcs when the words are laid out in their linear order, with all arcs above the words.

* PHRASE CHUNKING: find all non-recursive NPs and VPs in a sentence.

→ tag individual words with one of 3 tags:

- B (begin): word that starts a new target phrase
- I (inside): word that is part of target phrase but not the first word.
- O (other): word that is not part of target phrase.

- Evaluating chunking:

$$\text{Precision} = \frac{\# \text{ of correct chunks found}}{\text{Total \# of chunks found.}}$$

$$\text{Recall} = \frac{\# \text{ of correct chunks found}}{\text{Total \# of actual chunks.}}$$

→ F measure: harmonic mean of the two.

* Machine Translation

- Rule-Based: Hand-written transfer rules
 - Rules can be based on lexical / structural transfer
 - ⊕ Firm grip on complex translation phenomena
 - ⊖ often very labor-intensive → lack of robustness
- Statistical: Mainly word or phrase-based translations
 - translations are learned from actual data
 - ⊕ Translations are learned automatically
 - ⊖ Difficult to model complex translation phenomena.

→ Challenges:

- Sparsity: training models need a lot of data
 - genre & domain sensitive
- Hallucinations: ^{worse for language with rich morphology} cause unintentionally ~~biased~~ biased translations

- Automatic Evaluation: BLEU METRIC

- Modified n-gram precision with length penalty
- quick, inexpensive & language independent
- correlates highly with human evaluation
- bias against synonyms & inflectional variations.

- Word ~~Alignment~~ Alignment:

- Given parallel corpus, with word alignment → compute translation rules
- Given translation rules → compute probability of translation.
- Mapping alignment: given source word position i and target word position j ,

$$a: i \rightarrow j \Rightarrow a: \{1 \rightarrow 1, 2 \rightarrow 3, 3 \rightarrow 2, 4 \rightarrow 4\}$$

- Word-based translation model:

Initialize translation rules

Do until (convergence) {

Translate parallel corpus

Count translations

} Create translation rules

⇒ Does not take into account:

- Dropping words
- Inserting
- Word order
- one-to-many translation.

- Phrase-based translation model: extract phrase-pair consistent with word alignment

→ all words of the phrase-pair have to align to each other and no other word.

- Collect all phrase-pairs from the data

→ calculate probabilities: $t(e|f) = \frac{\text{Count}(e, f)}{\sum_i \text{Count}(e_i, f)}$

translation probability (pointing to $t(e|f)$)

target sentence (pointing to e)

source sentence (pointing to f)

- Evaluation:

- Automatic Metrics:

→ automatically optimize system performance towards metric.

- Goals: low cost, tunable, meaningful, consistent, correct
 - basic strategy: given machine translation output & human reference translation, compute the similarity between them.
 - Using precision & recall to compute harmonic $\frac{2}{1 + \frac{1}{\text{precision}} + \frac{1}{\text{recall}}}$ average ⇒ F-measure.
- FLAW: no penalty for re-ordering

- Word Error Rate : minimum number of editing steps to transform output to reference
 - Match: no cost
 - Substitution, insertion, deletion:
$$\text{Levenshtein distance (WER)} = \frac{\text{substitutions} + \text{insertions} + \text{deletions}}{\text{reference length}}$$
- BLEU : compute precision for n-gram of size 1 to 4
 - penalizes short translations \rightarrow brevity penalty
 - over entire test set, not individual sentences.
$$\text{BLEU} = \min \left(1, \frac{\text{output length}}{\text{reference length}} \right) \left(\prod_{i=1}^4 \text{precision}_i \right)^{1/4}$$
 - multiple references can be used
 - \rightarrow n-grams may match any of the references.
- Critique of Automatic Metrics: ignores relevance of words,
 - operates on local level \rightarrow do not consider overall grammaticality / sentence meaning
 - scores are meaningless \rightarrow test-specific
 - Human translators score low on BLEU

- PHRASE-BASED MODELS

- Advantages: many-to-many translation can handle non-compositional phrases
 - use of local context in translation
 - \rightarrow the more data, the longer phrases can be learned.
 - Phrase translation table: table with phrase translations & their probabilities.
 - \rightarrow will include lexical variations, morphological variations, include function words and noise.
 - Tuning: the process of finding the optimal weights for the linear decoding model
 - \rightarrow optimal weights are those which maximize translation performance on a small set of parallel sentences (tuning set)
 - MERT (Minimum Error Rate Training): decodes the whole training set and generates an n-best list \rightarrow model weights updated based on this decoder output
 - \rightarrow optimization process repeats until some convergence criterion is met
- transforming the source sentence into target sentence
 \rightarrow find best with highest probability.
- Methods of search space reduction:
 - Recombination: when two hypothesis paths lead to two matching hypotheses; if same number of foreign words are translated and same target words are in the output with different scores, drop worse hypothesis
 - Pruning: remove bad hypotheses early by putting ~~comparable~~ comparable hypothesis into stacks & limiting # of hypotheses in each stack
 - \rightarrow that have translated same # of input words.
- \Rightarrow Quadratic complexity

- Tokenization: solution for SMT being not as effective for language with rich morphology
→ break up the words to symmetrize source & target language.

* LEXICAL SEMANTICS

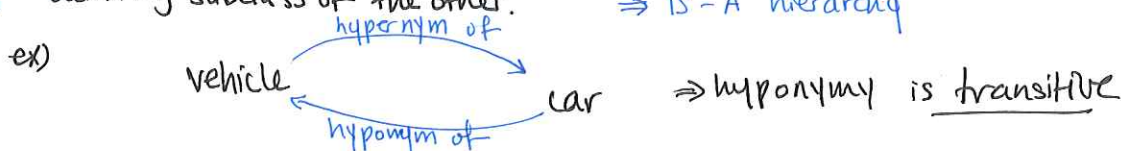
- Word sense: a discrete representation of an aspect of a word's meaning
- Homonyms: words that share a form but have unrelated, distinct meanings.
→ homographs (written same, read different), homophones (written different, read same)
- Polysemy: a polysemous word has related meanings
→ lots of types of polysemy are systematic. ex) school, hospital...
⇒ building ↔ organization

⇒ Do "zeugma" test to find out if a word has more than one sense.

ex) Does Lufthansa serve breakfast & San Jose?

- Synonyms: two lexemes are synonyms if they can be substituted for each other in all situations ⇒ they have the same propositional meaning
↔ Antonym

- Hyponymy: one sense is hyponym of another if the first sense is more specific, denoting subclass of the other. ⇒ IS-A hierarchy
↔ Hypernymy



* Similarity algorithms

thesaurus-based algorithms

- are words "nearby" in hypernym hierarchy?
- do they have similar glosses?

Distributional algorithms

- do words have similar distributional contexts?

$$\text{Sim}_{\text{path}} = \frac{1}{\text{pathlen}(c_1, c_2)}$$

- Path based similarity: two concepts are similar if they are near each other in the thesaurus hierarchy (concepts have path 1 to themselves)
→ problem: assume each link represents a uniform distance.

- Information content similarity:

- $P(c)$: the probability that a randomly selected word in a corpus is an instance of concept c → $P(\text{root}) = 1$

$$P(c) = \frac{\sum_{w \in \text{words}(c)} \text{count}(w)}{N}$$

words(c): set of all words that are children of node c

N → total number of words

- Information content (IC) = $-\log P(c)$

- Lowest Common Subsumer (LCS - most informative subsumer):

$\text{LCS}(c_1, c_2)$ = the most informative (lowest) node in the hierarchy subsuming both c_1 and c_2

$$\text{Sim}_{\text{resnik}}(c_1, c_2) = -\log P(\text{LCS}(c_1, c_2)) \rightarrow \text{the IC of the most informative subsumer of the two nodes.}$$

$$\text{Sim}_{\text{lin}}(c_1, c_2) = \frac{2 \log P(\text{LCS}(c_1, c_2))}{\log P(c_1) + \log P(c_2)} \propto \frac{\text{IC}(\text{common}(c_1, c_2))}{\text{IC}(\text{description}(c_1, c_2))}$$

→ the more differences b/w c_1 and c_2 , the less similar they are
amount of information needed to state the commonality
information needed to fully describe what A & B are.

$$\text{Sim}_{\text{elask}}(c_1, c_2) = \sum_{r, q \in \text{RELS}} \text{overlap}(\text{gloss}(r(c_1)), \text{gloss}(q(c_2)))$$

→ two concepts are similar if their glosses contain similar words

→ compute overlap for other relations (RELS) as well ⇒ glosses of hypernyms & hyponyms

- Distributional models of meaning (vector-space models of meaning)

→ offer much higher recall than hand-built thesauri (but lower precision)

• Term-document matrix: each cell is count of term t in a document d : $tf_{t,d}$

→ each document is a count vector in \mathbb{N}^V → two documents are similar if their vectors are similar

→ each word is a count vector in \mathbb{N}^D → two words are similar if their vectors are similar.

• Term-context matrix: use smaller contexts (paragraph, sentences) instead of entire documents

⇒ Instead of using raw counts: term-document matrix → TF-IDF (Term Frequency - Inverse Document Frequency)

Term-context matrix → Positive Pointwise Mutual Information (PPMI)

- Pointwise Mutual Information: do events X and Y co-occur more than if they were independent?

$$\text{PMI}(X, Y) = \log_2 \frac{P(X, Y)}{P(X)P(Y)}$$

⇒ For PPMI, replace all negative values with zero.

* Information Retrieval:

- Indexing: mapping between terms & documents.

→ "query" is a Boolean expression over vectors.

↳ For boolean retrieval, term-document matrix is binary ⇒ marks whether the term appears in the document or not.

⊕ Accurate with right strategies

⊕ Efficient for the computer

⊖ Results too many or none.

⊖ Users must know boolean logic

⊖ doesn't account for the fact that words have multiple meanings

• Indexing counts: entries are counts of occurrences of a term in a document.

- Document frequency: documents are most likely described well by rare terms that occur in them frequently

→ High Term Frequency (TF) is evidence of meaning

→ Low Document Frequency (DF) is evidence of importance.

⇒ High Inverse Document Frequency (IDF)

⇒ Term Weight = TF-IDF (product of TF & IDF)

$$w_{t,d} = tf_{t,d} \times \log \frac{N}{df_t} \leftarrow \# \text{ of documents in which "t" appears (normalized, inverted \& scaled)}$$

* TEXT SUMMARIZATION:

- simple baseline: take the first sentence

- 3 stages:
- 1) content selection: choose sentences to extract from the document
 - 2) information ordering: choose an order to place them in the summary
 - 3) sentence realization: clean up the sentences.
- baseline: use original of sentences

→ baseline: use order of appearance in doc

- CONTENT SELECTION:

• UNSUPERVISED CONTENT SELECTION: CHOOSE SENTENCES THAT HAVE SALIENT/INFORMATIVE WORDS

→ tf-idf: $W_i = tf_{ij} \times idf_j$

→ TOPIC SIGNATURE: CHOOSE A SMALLER SET OF SALIENT WORDS

• MUTUAL INFORMATION

• LOG-LIKELIHOOD RATIO (LLR):

$$W_i = \begin{cases} 1 & \text{if } -2 \log \lambda(W_i) > 10 \\ 0 & \text{otherwise} \end{cases}$$

• SUPERVISED CONTENT SELECTION: *hard to get*

GIVEN A LABELED TRAINING SET OF GOOD SUMMARIES FOR EACH DOCUMENT, ALIGN THE SENTENCES IN DOC WITH SENTENCES IN THE SUMMARY

← difficult

→ EXTRACT: POSITION OF SENTENCE, LENGTH OF SENTENCE, WORD INFORMATIVENESS, COHESION

→ TRAIN A BINARY CLASSIFIER (PUT SENTENCE IN SUMMARY?)

→ performance not better than unsupervised.

- ROUGE (Recall Oriented Understudy for Gisting Evaluation): based on BLEU INTRINSIC METRIC FOR AUTOMATICALLY EVALUATING SUMMARIES.

GIVEN DOCUMENT D , AUTOMATIC SUMMARY X ,

$$ROUGE_2 = \frac{\sum_{S \in \{RefSum\}} \sum_{\text{bigrams } i \in S} \min(\text{count}(i, X), \text{count}(i, S))}{\sum_{S \in \{RefSum\}} \sum_{\text{bigrams } i \in S} \text{count}(i, S)}$$

← bigrams

← human produced