

## 2.9 N-Gram Data Structures

- \_ N-Grams can be viewed as a special technique for conflation (stemming).
- \_ It is an unique data structure in information systems that ignores words and treats the input as a continuous data, optionally limiting its processing by interword symbols.
- \_ N-Grams are a fixed length consecutive series of “n” characters or fixed length overlapping symbol segments that define the searchable processing tokens.
- \_ These tokens have logical linkages to all the items in which the tokens are found.
- \_ Inversion lists, document vectors and other proprietary data structures are used to store the linkage data structure and are used in the search process.
- \_ Unlike stemming that generally tries to determine the stem of a word that represents the semantic meaning of the word, n-grams do not care about semantics.
  
- \_ Examples of bigrams, trigrams and pentagrams are given in Figure 4.7 for the word phrase “sea colony”.
- \_ For n-grams, with n greater than two, some systems allow interword symbols to be part of the n-gram set usually excluding the single character with interword symbol option.
- \_ The symbol # is used to represent the interword symbol which is anyone of a set of symbols (e.g., blank, period, semicolon, colon, etc.).

se ea co ol lo on ny	Bigrams (no interword symbols)
sea col olo lon ony	Trigrams (no interword symbols)
#se sea ea# #co col olo lon ony ny#	Trigrams (with interword symbol #)
#sea# #colo colon olony lony#	Pentagrams (with interword symbol #)

Figure 4.7 Bigrams, Trigrams and Pentagrams for “sea colony”

- \_ Each of the n-grams created becomes a separate processing token and are searchable.
- \_ It is possible that the same n-gram can be created multiple times from a single word.

- \_ The advantage of n-grams
  - \_ they place a finite limit on the number of searchable tokens.
  - \_ maximum number of unique n-grams can be generated.
  - \_ implementation techniques allow for fast processing on minimally sized machines
- \_ disadvantage
  - \_ false hits can occur under some architectures.
  - \_ increased size of inversion lists (or other data structures) that store the linkage data structure.
  - \_ In effect, use of n-grams expands the number of processing tokens by a significant factor.
  - \_ There is no semantic meaning in a particular n-gram since it is a fragment of processing token and may not represent a concept.
  - \_ Thus n-grams are a poor representation of concepts and their relationships.