

2.4 Automatic Indexing

- _ Advantages of human indexing
 - _ ability to determine concept abstraction
 - _ ability to judge the value of a concept.
- _ Disadvantages of human indexing
 - _ Cost
 - _ processing time
 - _ consistency

- _ *Automatic indexing*
 - _ Capability for the system to automatically determine the index terms to be assigned to an item.
 - _ More complex processing is required when emulate a human indexer and determine a limited number of index terms.
 - _ No additional indexing costs versus the salaries and benefits regularly paid to human indexers.
 - _ Requires only a few seconds or less of computer time based upon the size of the processor and the complexity of the algorithms to generate the index.
 - _ Consistency in the index term selection process as indexing is performed automatically by an algorithm.
 - _ Indexes from automated indexing fall into two classes:
 - _ unweighted.
 - _ weighted
- _ *Unweighted indexing system*
 - _ index term in a document and its word location(s) are kept in the searchable data structure.
 - _ Queries against unweighted systems are based upon Boolean logic and the items in the resultant Hit file are considered equal in value.
 - _ The last item presented in the file is the first item to be relevant to the user's information need.

- _ *Weighted indexing system*

- _ weight of the index term is based upon a function associated with the frequency of occurrence of the term in the item.
 - _ values for the index terms are normalized between zero and one.
 - _ The higher the weight, the more the term represents a concept discussed in the item.
 - _ query process uses the weights along with any weights assigned to terms in the query to determine a rank value.

- _ 2.4.1 Indexing by Term
- _ 2.4.2 Indexing by Concept
- _ 2.4.3 Multimedia Indexing

Indexing by Term

- _ There are two major techniques for creation of the index terms:
 - _ statistical
 - _ natural language
- _ Statistical techniques
 - _ based upon vector models and probabilistic models with a special case being Bayesian models.
 - _ calculation of weights in those models use statistical information such as the frequency of occurrence of words and their distributions in the searchable database.

- _ Natural language techniques
 - _ perform more complex parsing to define the final set of index concepts.
 - _ weighted systems are discussed as vectorized information systems.

- _ Vector model (Example. SMART system)
 - _ The system emphasizes weights as a foundation for information detection and stores these weights in a vector form.
 - _ Each vector represents a document
 - _ each position in a vector represents a different unique word (processing token) in the database.

- _ The value assigned to each position is the weight of that term in the document.
- _ A value of zero indicates that the word was not in the document.
- _ Queries can be translated into the vector form.
- _ Search is accomplished by calculating the distance between the query vector and the document vectors.
- _ Probabilistic model.
 - _ Bayesian approach is the most successful model in this area.
 - _ It is based upon the theories of evidential reasoning (drawing conclusions from evidence).
 - _ The Bayesian approach could be applied as part of index term weighting by calculating the relationship between an item and a specific query.
 - _ A Bayesian network is a directed acyclic graph
 - _ Each node represents a random variable.
 - _ Arcs between the nodes represent a probabilistic dependence between the node and its parents.
- _ Figure shows the basic weighting approach for index terms or associations between query terms and index terms.

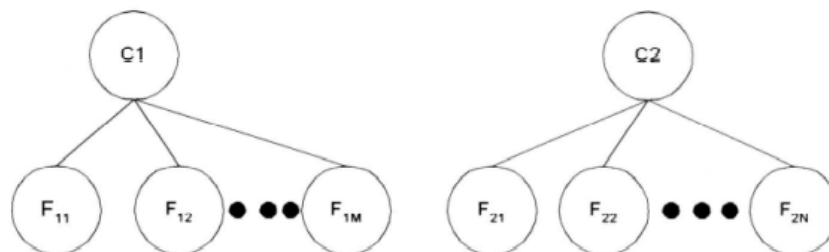


Figure 3.3 Two-level Bayesian network

- _ The nodes C1 and C2 represent “the item contains concept C_i ”
- _ F nodes represent “the item has feature F_{ij} (e.g., words)”
- _ The network interpreted as C representing concepts in a query
- _ F representing concepts in an item.
- _ The goal is to calculate the probability of C_i given F_{ij} .
- _ To perform that calculation two sets of probabilities are needed:

- _ The prior probability $P(C_i)$ that an item is relevant to concept C .
- _ The conditional probability $P(F_{ij}/C_i)$ that the features F_{ij} where $j=1,m$ are present in an item given that the item contains topic C_i .
- _ The automatic indexing task is to calculate the posterior probability $P(C_i/F_{i1}, \dots, F_{im})$ the probability that the item contains concept C_i given the presence of features F_{ij}

- _ The Bayes inference formula that is used is:

$$P(C_i/F_{i1}, \dots, F_{im}) = P(C_i) P(F_{i1}, \dots, F_{im}/C_i) P(F_{i1}, \dots, F_{im}).$$

- _ If the goal is to provide ranking as the result of a search by the posteriors, the Bayes rule can be simplified to a linear decision rule:

$$g(C_i/F_{i1}, \dots, F_{im}) = \sum_k I(F_{ik}) w(F_{ik}, C_i)$$

- _ where $I(F_{ik})$ is an indicator variable that equals 1 only if F_{ik} is present in the item (equals zero otherwise) and w is a coefficient corresponding to a specific feature/concept pair.
- _ function g is the sum of the weights of the features
- _ w is weight corresponding to each feature (index term)
- _ w produces a ranking in decreasing order that is equivalent to the order produced by the posterior probabilities.

- _ DR-LINK (Document Retrieval through Linguistic Knowledge) system
 - _ define indexes to items via natural language processing.
 - _ processes items at the morphological, lexical, semantic, syntactic, and discourse levels.
 - _ Each level uses information from the previous level to perform its additional analysis.
 - _ The discourse level is abstracting information beyond the sentence level and can determine abstract concepts.
 - _ This allows the indexing to include specific term as well as abstract concepts such as time.

Indexing by Concept

- _ Indexing by term treats each occurrence as a different index.

- _ Then uses thesauri or other query expansion techniques to expand a query to find the different ways the same thing has been represented.
- _ *basis for concept indexing*
 - _ There are many ways to express the same idea and increased retrieval performance comes from using a single representation.
 - _ Concept indexing determines a related set of concepts based upon a test set of terms
 - _ uses those concepts for indexing all items
 - _ Latent Semantic Indexing
 - _ Indexing the latent semantic information in items.
- _ *MatchPlus system.*
 - _ Example for concept indexing.
 - _ neural networks facilitates machine learning of concept/word relationships.
 - _ goal is to determine word relationships (e.g., synonyms) and the strength of these relationships and use that information in generating context vectors, from the corpus of items.
 - _ Two neural networks are used.
 - _ One neural network learning algorithm generates stem context vectors that are sensitive to similarity of use.
 - _ another one performs query modification based upon user feedback.
 - _ *context vectors*
 - _ Word stems, items and queries that are represented by high dimensional (atleast 300 dimensions) vectors.
 - _ Each dimension in a vector could be viewed as an abstract concept class.

Multimedia Indexing

- _ The automated indexing takes place in multiple passes of the information versus just a direct conversion to the indexing structure.

- _ The first pass in most cases is a conversion from the analog input mode into a digital structure.
- _ Then algorithms are applied to the digital structure to extract the unit of processing of the different modalities that will be used to represent the item.
- _ In an abstract sense this could be considered the location of a processing token in the modality.
- _ This unit will then undergo the final processing that will extract the searchable features that represent the unit.

- _ Indexing video or images can be accomplished
 - _ at the raw data level (e.g., the aggregation of raw pixels),
 - _ the feature level distinguishing primitive attributes such as color and luminance
 - _ at the semantic level where meaningful objects are recognized (e.g., an airplane in the image/video frame).

- _ An example is processing of video.
- _ The system (e.g., Virage) will periodically collect a frame of video input for processing.
- _ It might compare that frame to the last frame captured to determine the differences between the frames.
- _ If the difference is below a threshold it will discard the frame.
- _ For a frame requiring processing, it will define a vector that represents the different features associated with that frame.
- _ Each dimension of the vector represents a different feature level aspect of the frame.
- _ The vector then becomes the unit of processing in the search system.
- _ This is similar to processing an image.
- _ Semantic level indexing requires pattern recognition of objects within the images.

- _ analog audio input
 - _ system will convert the audio to digital format and determine the phonemes associated with the utterances.

- _ The phonemes will be used as input to a Hidden Markov Search model that will determine with a confidence level the words that were spoken.
 - _ A single phoneme can be divided into four states for the Markov model.
 - _ It is the textual words associated with the audio that becomes the searchable structure.

- _ In addition to storing the extracted index searchable data, a multimedia item also needs to store some mechanism to correlate the different modalities during search.
- _ There are two main mechanisms that are used
 - _ Positional
 - _ Temporal
- _ Positional
 - _ is used when the modalities are scattered in a linear sequential composition.
 - _ For example a document that has images or audio inserted, can be considered a linear structure and the only relationship between the modalities will be the just a position of each modality.
- _ Temporal
 - _ based upon time because the modalities are executing concurrently.
 - _ The typical video source off television is inherently a multimedia source.
 - _ It contains video, audio, and potentially closed captioning.

- _ Synchronized Multimedia Integration Language (SMIL).
 - _ creation of multimedia presentations are becoming more common using the Synchronized Multimedia Integration Language (SMIL).
 - _ It is a mark-up language designed to support multimedia presentations that integrate text (e.g., from slides or free running text) with audio, images and video.
 - _ time is the mechanism that is used to synchronize the different modalities.
 - _ indexing must include a time-offset parameter versus a physical displacement.