

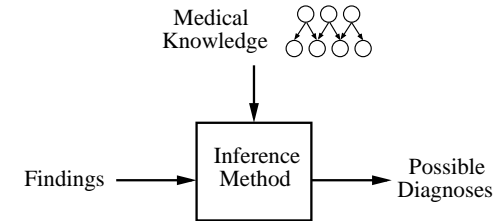
LECTURE 25:  
CLINICAL AND BIOINFORMATICS APPLICATIONS

Sam Roweis

April 2, 2004

MEDICAL DIAGNOSIS

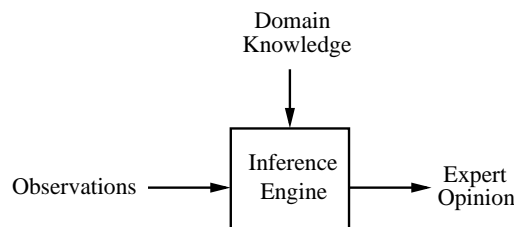
---



- In medical diagnosis the observations are clinical findings (what's wrong with this person), the domain knowledge represents which diseases or ailments have which symptoms, as well as which diseases are most likely for certain types of patients.
- The expert opinion takes the form of possible diagnosis (what ailment is most likely to be causing their problems).

GENERAL MOTIVATION FOR EXPERT SYSTEMS

---



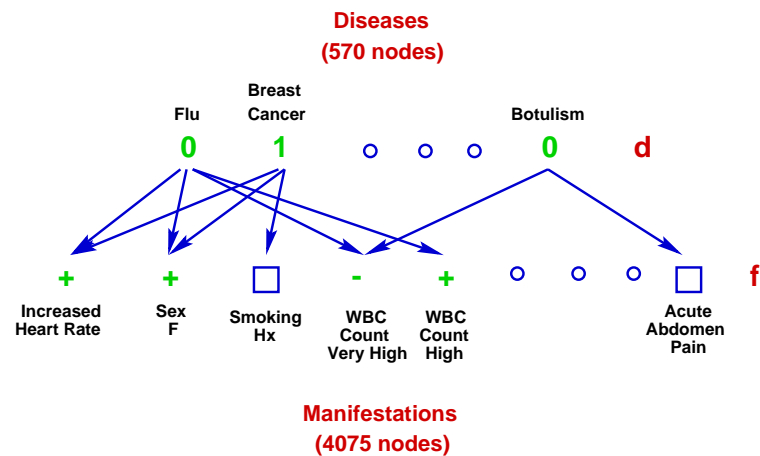
- Expert Systems attempt to combine domain knowledge with noisy observations and use a rational inference engine (often probabilistic) to come up with a conclusion or opinion.
- The two main problems in expert systems are how to encode the domain knowledge and how to perform inference efficiently.

QUICK MEDICAL REFERENCE (QMR-DT)

---

- Quick Medical Reference, Decision Theoretic (QMR-DT) is a very large graphical model based on expert knowledge acquired from medical doctors and clinical records in hospitals.
- There are 570 diseases and 4075 manifestations, which include symptoms, demographic data about the patient, medical history, and results of laboratory tests.
- We represent these using binary random variables  $d_k$  and  $f_i$ , encoding all non-binary manifestations (e.g. continuous values or categorical findings) with one-hot or range values.
- The domain knowledge was *not* learned from data directly using maximum likelihood, etc. Instead it was captured from the historical medical literature and from expert opinions and encoded into the graphical model by hand.

## QUICK MEDICAL REFERENCE (QMR-DT)



## INFERENCE IS THE KEY

- The full posterior is huge (exponential in the number of diseases), so we can only ever hope to compute its marginals.
- Even just to compute the likelihood requires a large amount of work because we have to sum over all possible disease configurations.

$d$ : disease configuration

$f$ : findings

$$P(d | f) = \frac{P(d) P(f | d)}{\sum_{d'} P(d') P(f | d')}$$

prior
conditional
joint

posterior
probability of evidence

## QUICK MEDICAL REFERENCE (QMR-DT)

- The graphical model asserts that manifestations are conditionally independent given the diseases:

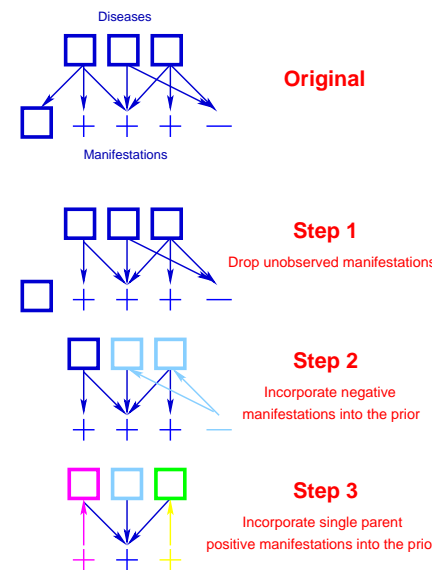
$$p(\mathbf{d}, \mathbf{f}) = \left[ \prod_k p(d_k) \right] \left[ \prod_i p(f_i | d) \right]$$

- The conditional model for activation of the manifestations given the diseases is noisy-OR:

$$\log p(f_i = 0 | \mathbf{d}) = w_{i0} + \sum_k w_{ik} d_k$$

- Most of the time very few diseases are active (less than 9), and zero or one diseases account for 72% of the mass under the disease prior.
- Also, usually between a few and a hundred manifestations are observed out of the 4075 possibilities.
- The noisy-OR weights are also very sparse: only 2% are nonzero.

## INFERENCE IN 2-LAYER BINARY NOISY-OR NETWORKS



But there is a trick...

- The Quickscore Algorithm (Heckerman 1989) computes  $P(f)$  in time exponential in number of positive findings with multiple parents.
- The trick is that negative findings and positive findings with only one parent can be absorbed into the prior.
- Still, with 100 observations we would still have to sum over  $2^{100}$  configurations.

## APPROXIMATE INFERENCE

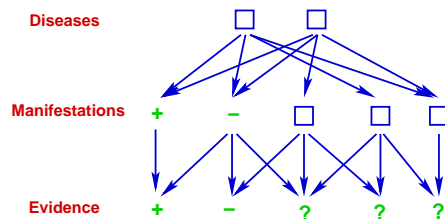
- Even with the quickscore trick, exact inference is often intractable in networks as large as QMR-DT.
- So practitioners resort to approximate inference methods which attempt to *estimate* the marginals  $p(d_k|\mathbf{f})$  rather than computing them exactly.
- This is a large and complex area of research, but essential to making QMR a practical diagnosis system.

## DISCRETE SEQUENCES IN COMPUTATIONAL BIOLOGY

- There has recently been a great interest in applying probabilistic models to analyzing discrete sequence data in molecular and computational biology.
- There are two major sources of such data:
  - amino acid sequences for protein analysis
  - base-pair sequences for genetic analysis
- The sequences are sometimes annotated by other labels, e.g. species, mutation/disease type, gender, race, etc.
- Lots of interesting applications:
  - whole genome shotgun sequence fragment assembly
  - multiple alignment of conserved sequences
  - splice site detection
  - inferring phylogenetic trees

## QMR-DT OBSERVATION PROCESS

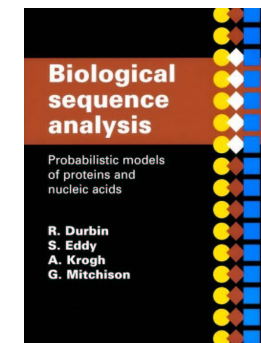
- Other issues: there should be a distinction between *unobserved* manifestations and *observed negative* manifestations.
- Observation is not independent of result: doctor's do the tests they expect will give them important info.
- Modeling this observation process is key to using QMR in practice (see work of Quaid Morris).



## MAIN TOOL: HIDDEN MARKOV MODELS

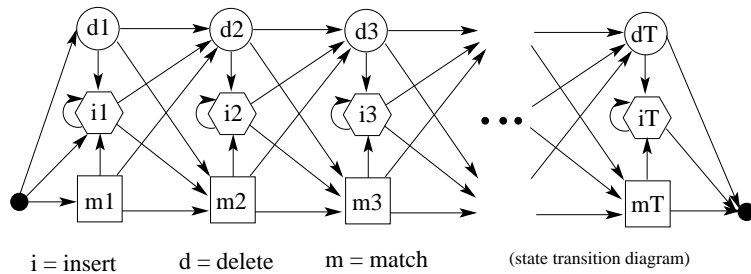
- HMMs and related models (e.g. profile HMMs) have been the major tool used in biological sequence analysis and alignment.
- The basic dynamic programming algorithms can be improved in special cases to make them more efficient in time or memory.

See the excellent book by Durbin, Eddy, Krogh, Mitchison for lots of practical details on applications and implementations.





## PROFILE (STRING-EDIT) HMMs



- A “profile HMM” or “string-edit” HMM is used for probabilistically matching an observed input string to a stored template pattern with possible insertions and deletions.
- Three kinds of states: match, insert, delete.
  - $m_j$  – use position  $j$  in the template to match an observed symbol
  - $i_j$  – insert extra symbol(s) observations after template position  $j$
  - $d_j$  – delete (skip) template position  $j$