# Hugin & Weka for Learning from Data

Using Bayesian Networks

Ashwin Patthi

University of South Carolina

# Learning from Data?

Predicting Probability distribution of unknown variables from data/cases

Resulting learning model can be used as classifier

Knowledge discovery

# Classification Task

Classifying variable y = x0 called the class variable given set of variables X=x1,x2…xk called attribute variables

Classifier is learned from a dataset D consisting of samples over (x,y) on network Bx over probability distribution U

$$h: X \longrightarrow y$$

# Learning using Hugin

Uses EM (Estimation Maximization) Algorithm (Batch learning)

This is used only when structure is available

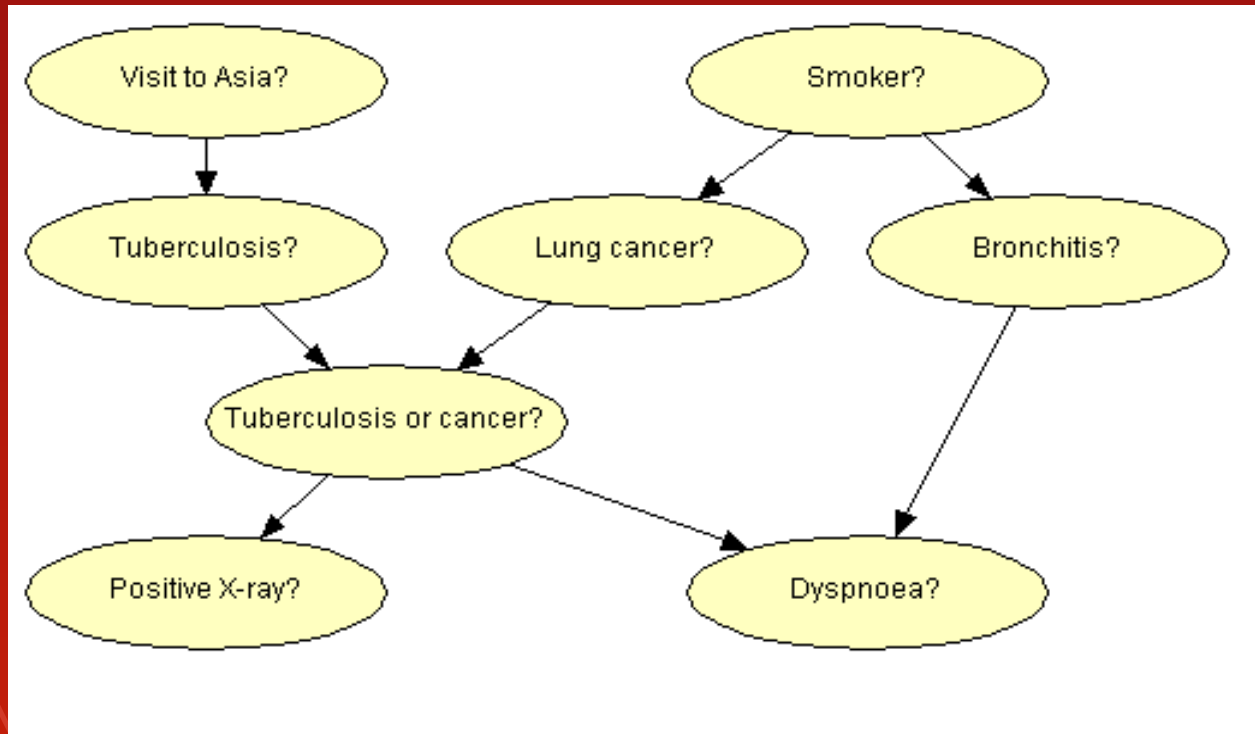Experience table must be provided for nodes whose conditional probabilities are to inferred

# EM Algorithm

Performs number of iteration on cases

Computes log-likelihood and attempts to maximize it

Stops when two successive log-likelihood is less than tolerance
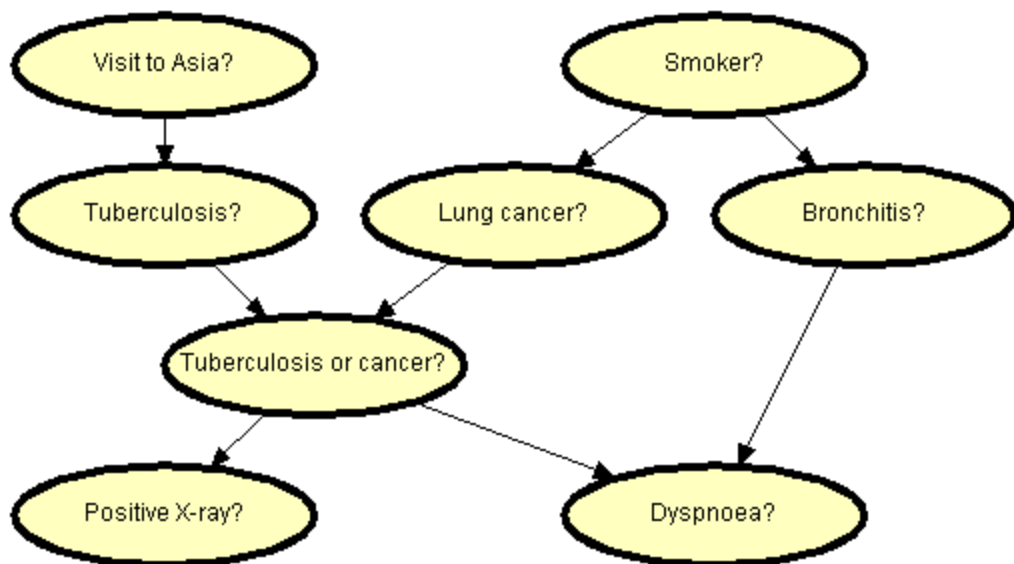
# Learning on Chest Clinic data

# Assumptions

No prior knowledge on any distribution

Set probability distribution to 1 except for "Tuberculosis or Cancer"

Also experience value to 10 (or some low value)

# Data set (asia.data)

First line is header

Each record is a case

N/A = Not available

```
asia.dat
 1  E,T,L,S,A,D,B,X
 2  no,no,no,yes,no,yes,yes,no
 3  no,no,no,yes,no,no,yes,no
 4  no,N/A,no,no,no,yes,yes,no
 5  no,no,no,no,no,no,no,no
 6  no,no,no,no,no,yes,yes,no
 7  no,no,no,yes,no,no,yes,no
 8  no,no,no,no,no,no,no,no
 9  no,no,no,yes,no,yes,no,no
10  no,no,no,no,no,no,no,no
11  no,no,no,yes,N/A,yes,yes,no
12  no,no,no,yes,N/A,yes,yes,no
13  no,no,no,yes,no,no,no,no
14  no,no,no,yes,no,no,no,no
15  yes,no,yes,yes,no,yes,yes,yes
16  N/A,no,no,N/A,no,yes,yes,N/A
17  no,no,no,no,no,yes,yes,no
18  no,no,N/A,yes,no,no,no,no
19  no,no,no,yes,no,no,no,no
20  no,no,no,no,no,N/A,no,no
21  no,no,N/A,N/A,no,yes,yes,no
22  no,no,no,no,no,no,no,no
23  yes,yes,no,yes,no,no,yes,yes
24  yes,no,yes,yes,no,no,no,yes
25  no,no,no,no,no,yes,yes,N/A
```

# Run the Learning Algo

# Resulting Marginal Probabilities

# Classifier

The resulting model with computed conditional probabilities can be used as classifier to predict the new unknown conditional probabilities
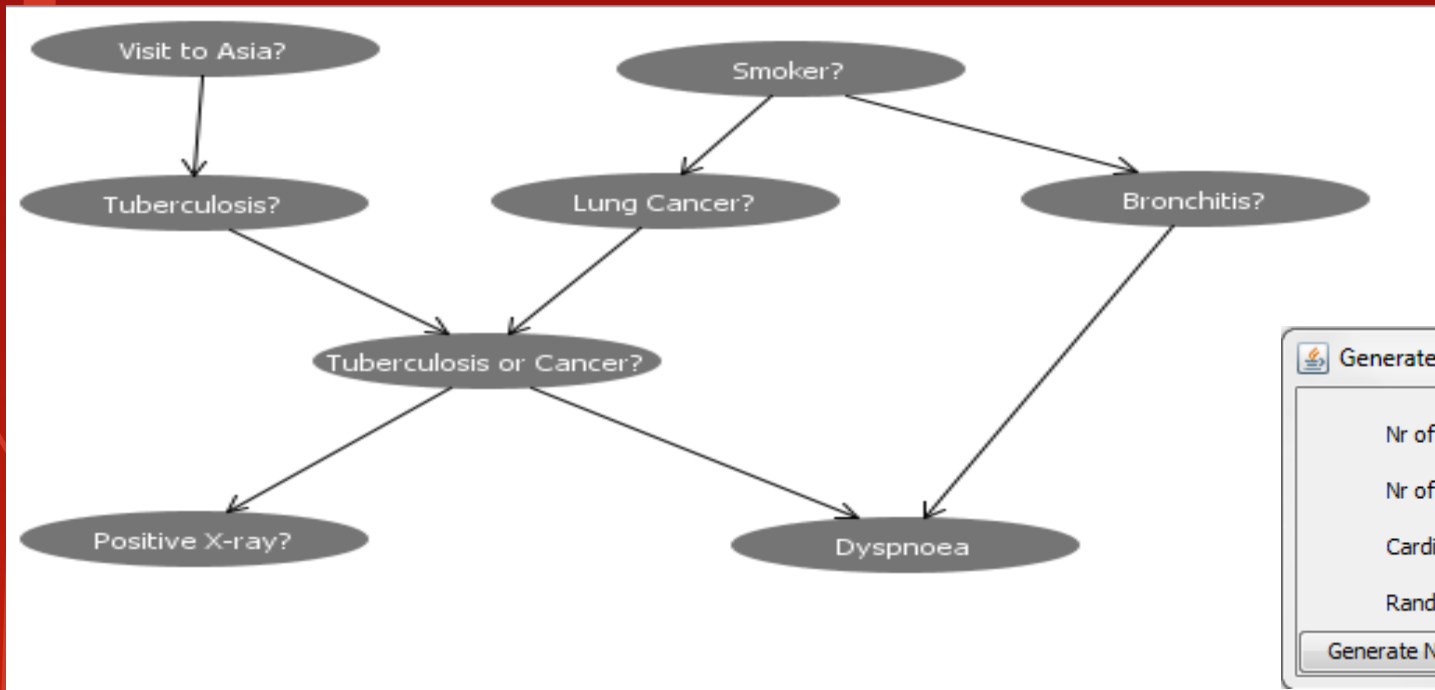
# Learning using WEKA

Collection of Machine Learning Algorithms for data mining tasks

Contains tools for data pre-processing, classification, clustering, visualization etc.

Bayesian Network classifier and editor is one of them

# Creating Bayesian Network

# Basic Assumptions

All variables are discrete finite variables

If continuous can convert to discrete using class filters.attribute.Discretize

No instances have missing values. If found can be filled by attribute.ReplaceMissingValues

# ✦ Results

# Pros & Cons of Hugin

**Pros:**

Can generate data with missing values based on available network

Can predict prob. distribution for all variables

**Cons:**

No feature for handling missing values

No feature to predict class of all the cases at once

# Pros and Cons of WEKA

Pros:

Feature to handle missing values

Can predict class of unknown data sets all at once

Cons:

Can't generate data with missing values

Can't predict probability distribution

# Useful Links & References

1. Bouckaert, Remco R. Bayesian network classifiers in weka. Department of Computer Science, University of Waikato, 2004

2. Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, Ian H. Witten (2009); The WEKA Data Mining Software: An Update; SIGKDD Explorations, Volume 11, Issue 1

3. More information on Hugin tool

4. More information on WEKA tool

Questions ?