# CS 4824/ECE 4424: Graphical Models II

**Acknowledgement**:

Many of these slides are derived from Tom Mitchell, Pascal Poupart, Pieter Abbeel, Eric Eaton, Carlos Guestrin, William Cohen, and Andrew Moore.

# Bayesian network recap

- A Bayes network represents the joint probability distribution over a collection of random variables

- A Bayes network is a directed acyclic graph and a set of CPD's
  - Each node denotes a random variable
  - Edges denote dependencies
  - CPD for each node $X_i$ defines $P(X_i | Pa(X_i))$
  - The joint distribution over all variables is defined as

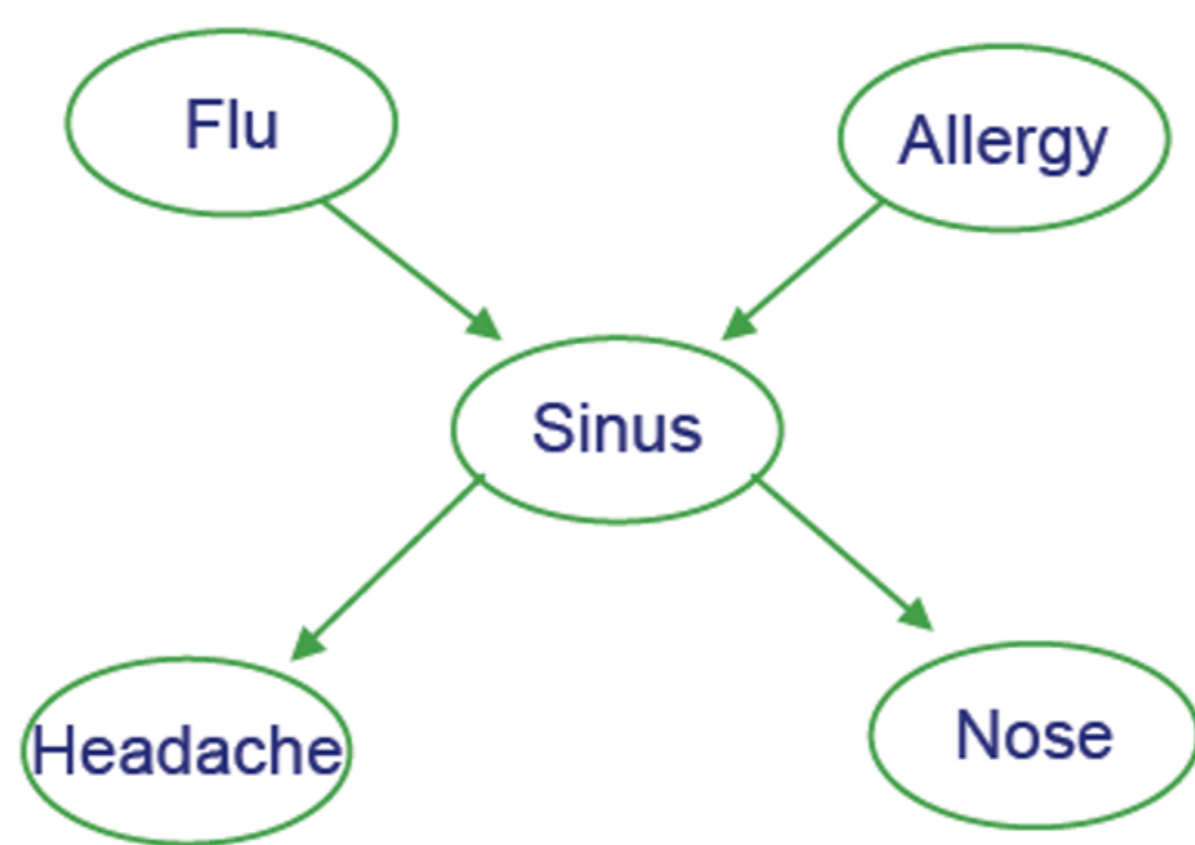$$P(X_1 \ldots X_n) = \prod_i P(X_i | Pa(X_i))$$

# Bayesian network recap

◦ Bayes nets are convenient representation for encoding dependencies / conditional independence

◦ BN = Graph plus parameters of CPD's
  ◦ Defines joint distribution over variables
  ◦ Can calculate everything else from that
  ◦ Though inference may be intractable

# Inference in Bayes Nets

- In general, intractable (NP-complete)
- For certain cases, tractable
  - Assigning probability to fully observed set of variables
  - Or if just one variable unobserved
  - Or for singly connected graphs (ie., no undirected loops)
    - Belief propagation
- Sometimes use Monte Carlo methods
  - Generate many samples according to the Bayes Net distribution, then count up the results
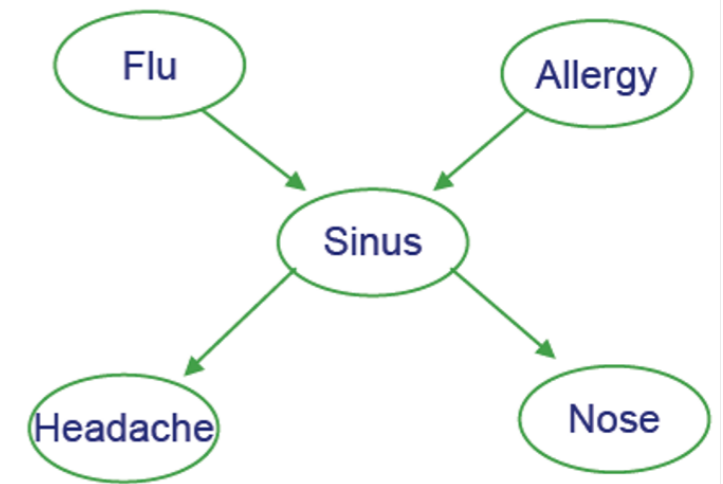
# Example

◦ Bird flu and Allegies both cause Sinus problems
◦ Sinus problems cause Headaches and runny Nose
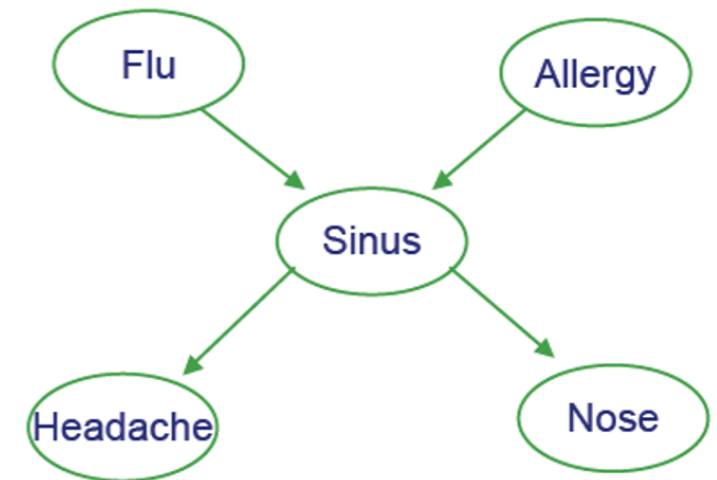
# Prob. of joint assignment: easy



- Suppose we are interested in joint assignment <F=f,A=a,S=s,H=h,N=n>

- What is P(f,a,s,h,n)?

let's use P(a,b) as shorthand for P(A=a, B=b)
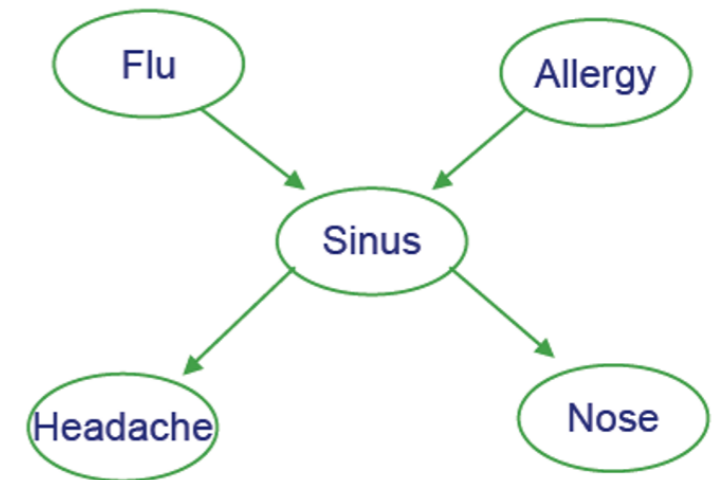
# Prob. of marginals: not so easy

○ How do we calculate P(N=n)?



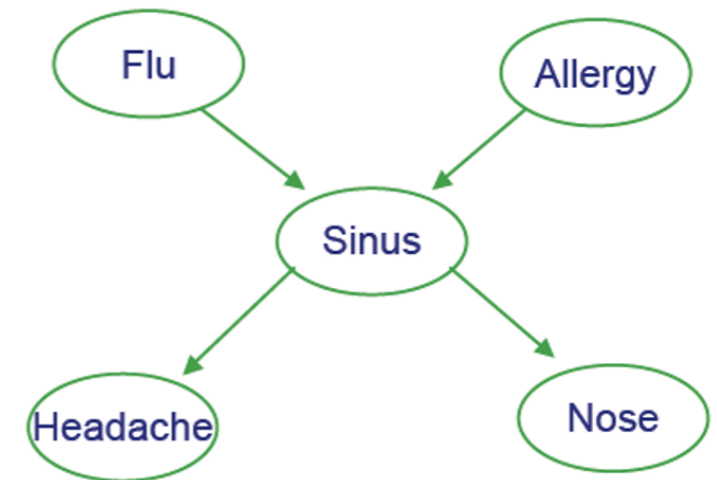let's use P(a,b) as shorthand for P(A=a, B=b)

# Generating a sample from joint distribution: easy

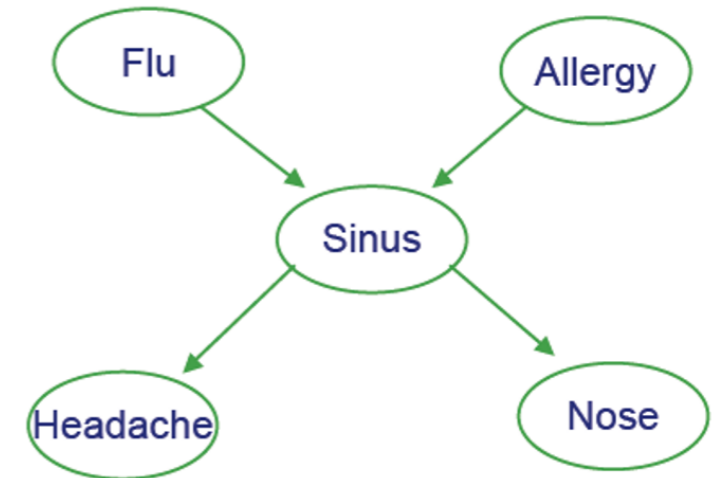○ How can we generate random samples drawn according to P(F,A,S,H,N)?

# Generating a sample from joint distribution: easy

- How can we generate random samples drawn according to P(F,A,S,H,N)?

- random sample of F according to $P(F=1) = \theta_{F=1}$ :
  - draw a value of r uniformly from [0,1]
  - if $r<\theta$ then output F=1, else F=0

- **Solution**:
  - draw a random value f for F, using its CPD
  - then draw values for A, for S|A,F, for H|S, for N|S

# Generating a sample from joint distribution: easy

◦ Note we can estimate marginals like $P(N=n)$ by generating many samples from joint distribution, then count the fraction of samples for which $N=n$



◦ Similarly, for anything else we care about $P(F=1|H=1, N=0)$
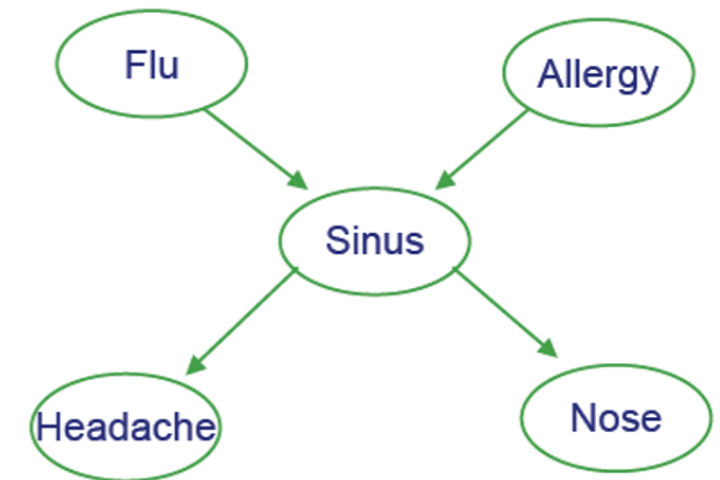
   ◦ weak but general method

# Learning of Bayes Nets

- Several types of of learning problems
  - Variable values may be fully observed / partly unobserved

- Easy case: learn parameters when data is *fully observed*

- Interesting case: graph *known*, data *partly known*

# Learning CPTs from Fully Observed Data

- Example: Consider learning the parameter
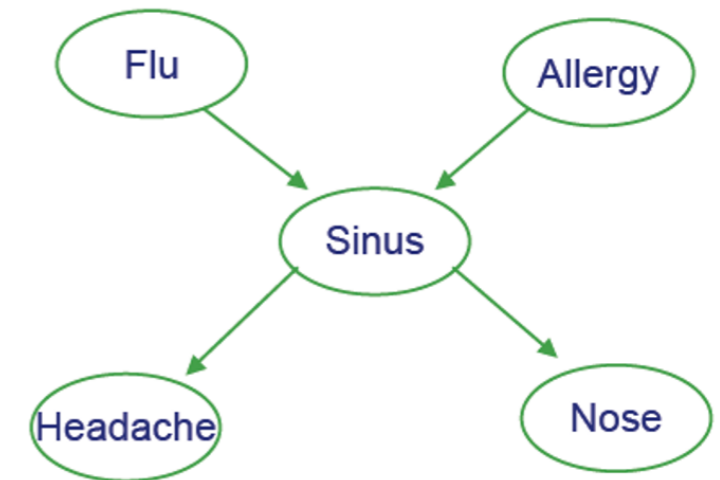
$$\theta_{s|ij} \equiv P(S = 1 | F = i, A = j)$$

- Max Likelihood Estimate is

$$\theta_{s|ij} = \frac{\sum_{k=1}^{K} \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^{K} \delta(f_k = i, a_k = j)}$$

# Learning CPTs from Fully Observed Data

◦ Example: Consider learning  the parameter

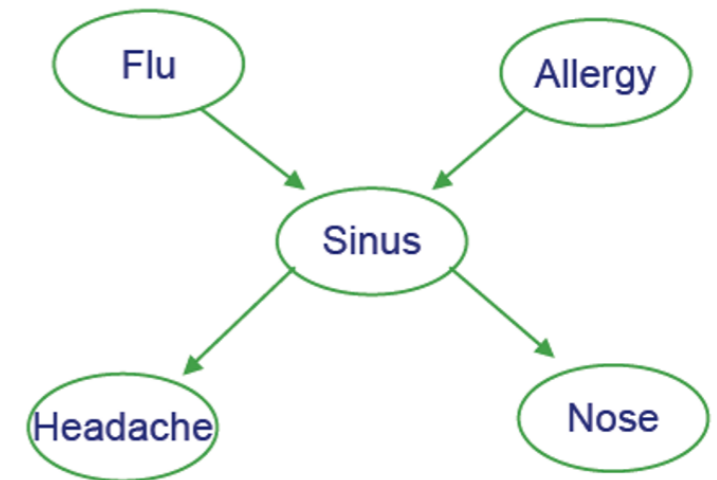$$\theta_{s|ij} \equiv P(S = 1 | F = i, A = j)$$



◦ Max Likelihood Estimate is

$$\theta_{s|ij} = \frac{\sum_{k=1}^{K} \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^{K} \delta(f_k = i, a_k = j)}$$

$k^{th}$ training  example

$\delta(x) = 1$ if x=true,
$\quad = 0$ if x=false

# MLE estimate of $\theta_{s|ij}$ from fully observed data
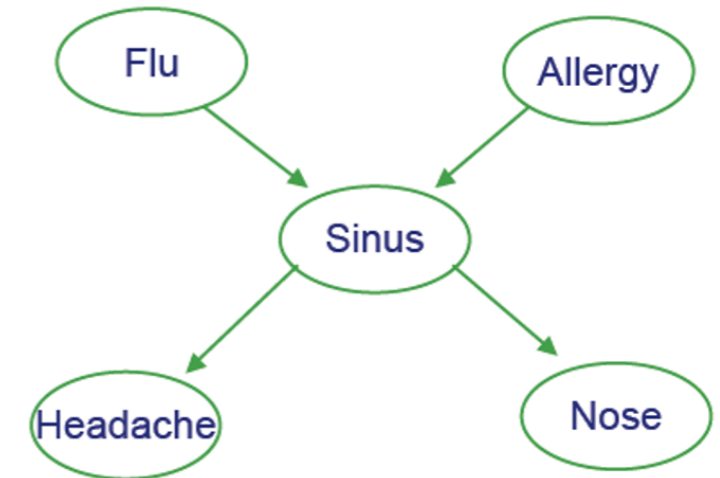


◦ Maximum likelihood estimate


◦ Our case

# MLE estimate of $\theta_{s|ij}$ from fully observed data

○ Maximum likelihood estimate

$$\theta \leftarrow \arg\max_{\theta} \log P(data|\theta)$$



○ Our case

$$P(data|\theta) = \prod_{k=1}^{K} P(f_k, a_k, s_k, h_k, n_k)$$

$$P(data|\theta) = \prod_{k=1}^{K} P(f_k)P(a_k)P(s_k|f_k a_k)P(h_k|s_k)P(n_k|s_k)$$

$$\log P(data|\theta) = \sum_{k=1}^{K} \log P(f_k) + \log P(a_k) + \log P(s_k|f_k a_k) + \log P(h_k|s_k) + \log P(n_k|s_k)$$

$$\frac{\partial \log P(data|\theta)}{\partial \theta_{s|ij}} = \sum_{k=1}^{K} \frac{\partial \log P(s_k|f_k a_k)}{\partial \theta_{s|ij}}$$

$$\theta_{s|ij} = \frac{\sum_{k=1}^{K} \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^{K} \delta(f_k = i, a_k = j)}$$

# Estimate $\theta$ from partly observed data

- What if FAHN observed, but not S?

- Can't calculate MLE

$$\theta \leftarrow \arg\max_{\theta} \log \prod_k P(f_k, a_k, s_k, h_k, n_k | \theta)$$

- Let X be all *observed* variable values (over all examples)
- Let Z be all unobserved variable values
- Can't calculate MLE

$$\theta \leftarrow \arg\max_{\theta} \log P(X, Z | \theta)$$

- What to do?