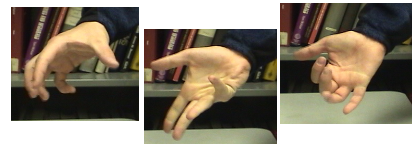
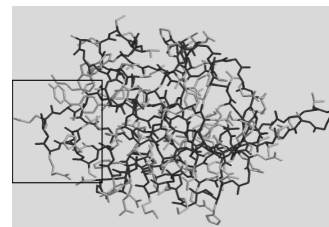
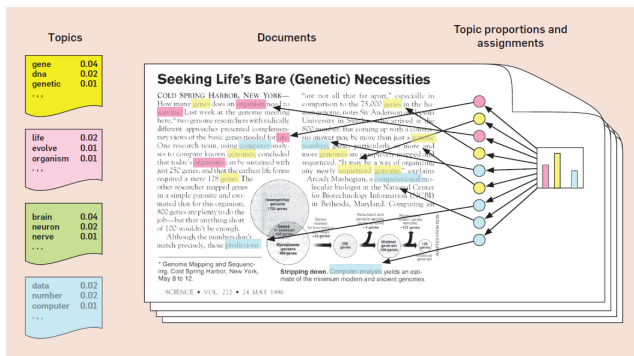


# Bayesian Networks— Representation

CSE 446: Machine Learning  
Emily Fox  
University of Washington  
March 6, 2017

©2017 Emily Fox

## Learning from structured data

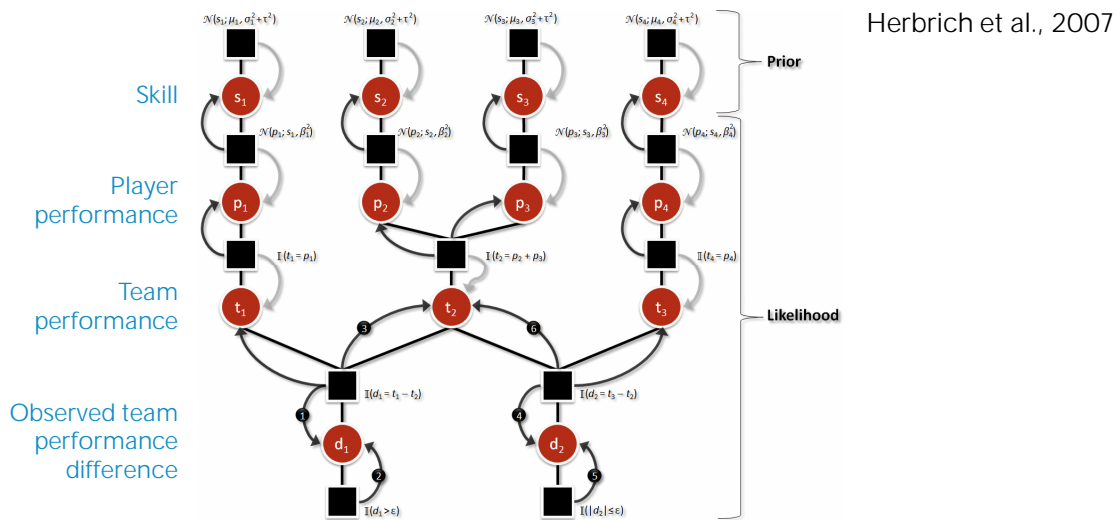


2

©2017 Emily Fox

CSE 446: Machine Learning

# TrueSkill: A Bayesian Skill Rating System

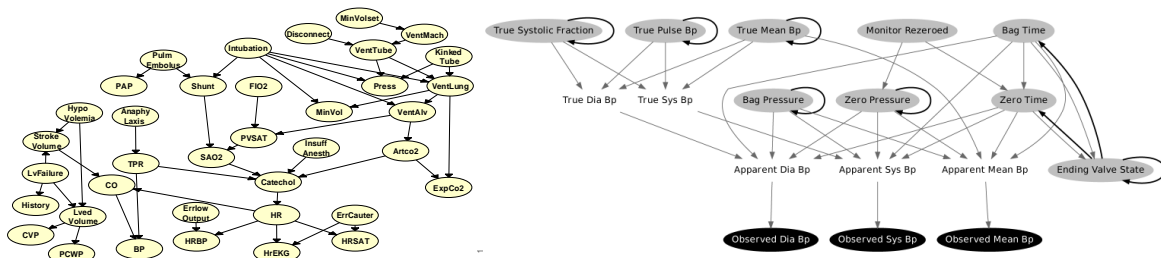


3

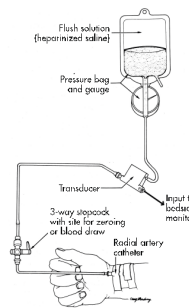
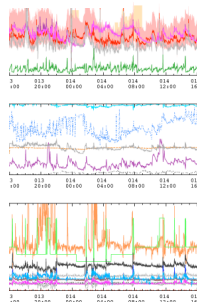
©2017 Emily Fox

CSF 446: Machine Learning

# ICU Monitoring



Beinlich et al., 1989



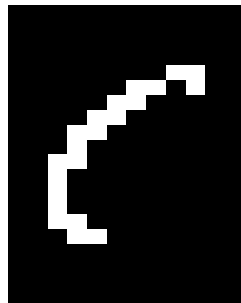
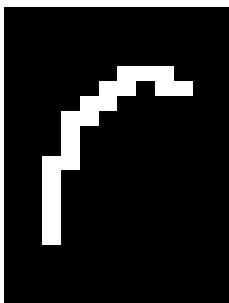
4

©2017 Emily Fox

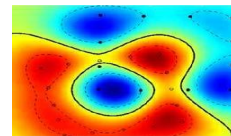
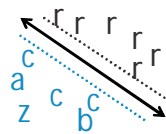
CSF 446: Machine Learning

## Digging in: Learning with and without context/structure

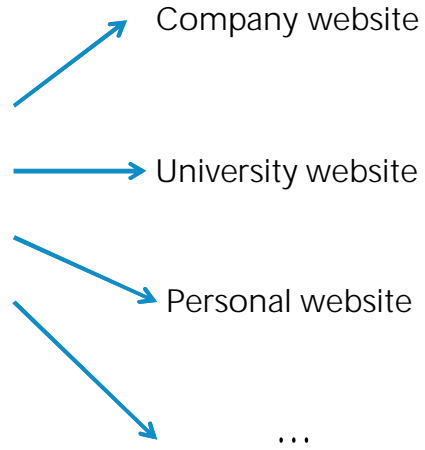
### Without context: Handwriting recognition



Character recognition,  
e.g., kernel SVMs



# Without context: Webpage classification

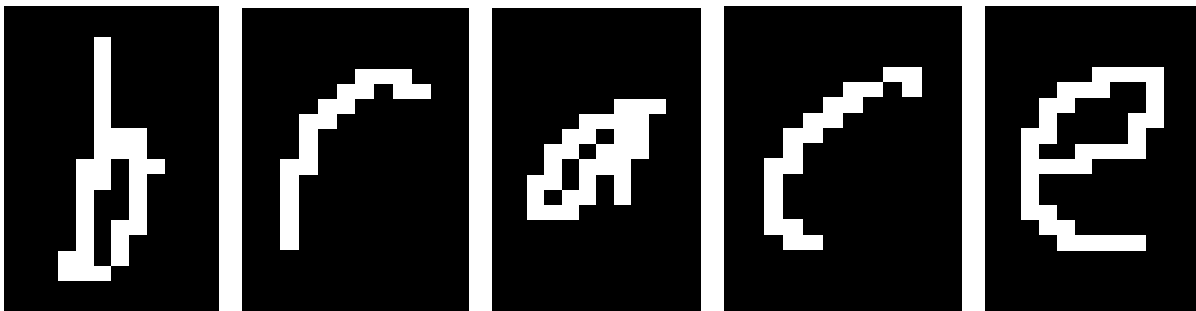


7

©2017 Emily Fox

CSF 446: Machine Learning

# With context: Handwriting recognition

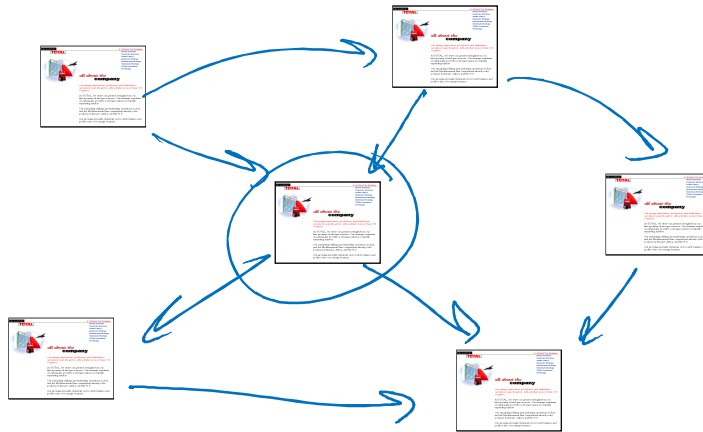


"c" more likely to come after "a" than another "a"

©2017 Emily Fox

CSF 446: Machine Learning

## With context: Webpage classification



*Company pages  
tend to  
point to  
each other*

9

©2017 Emily Fox

CSF 446: Machine Learning

Modeling structured relationships  
via Bayesian networks

## Today – Bayesian networks

- Provided a huge advancement in AI/ML
- Generalizes naïve Bayes and logistic regression
- ★ • Compact representation for exponentially-large probability distributions
- Exploit conditional independencies

11

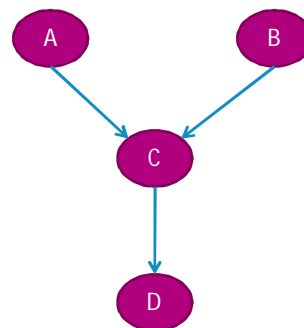
©2017 Emily Fox

CSF 446: Machine Learning

## Bayesian network representation

Compact representation of a probability distribution.

Directed Acyclic Graph



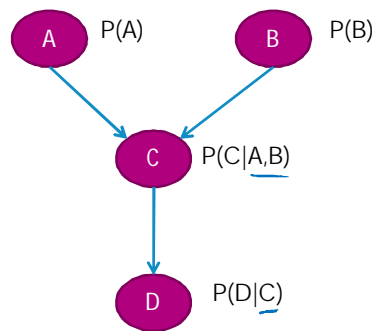
**Vertices:** Random Variables  
**Edges:** Conditional dependencies  
 "probabilistic relationships"

12

©2017 Emily Fox

CSF 446: Machine Learning

# Bayesian network probability factorization



One **CPT** (conditional probability table) for each variable

**P(variable | parents of variable)**

implies the factorization:

$$P(\mathbf{X}) = \prod_{j=1}^d P(\mathbf{X}[j] \mid \text{parents}(\mathbf{X}[j]))$$

↑
↑  
node
parents

joint:

$$P(A,B,C,D) = P(A) P(B) P(C|A,B) P(D|C)$$

13

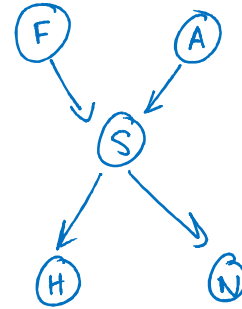
©2017 Emily Fox

CSF 446: Machine Learning

What a Bayesian network represents (in detail) and what does it buy you?

## Causal structure

- Suppose we know the following:
  - The **flu** causes **sinus** inflammation
  - **Allergies** cause **sinus** inflammation
  - **Sinus** inflammation causes a **runny** nose
  - **Sinus** inflammation causes **headaches**
- How are these connected?



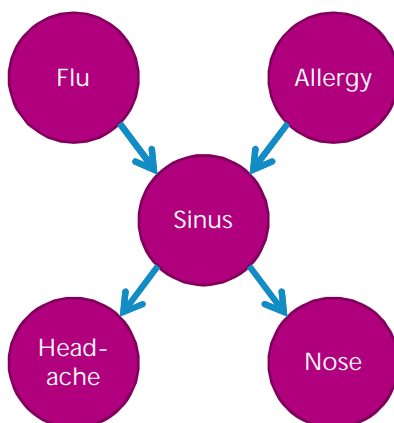
*(Not a true  
"causal model")*

15

©2017 Emily Fox

CSF 446: Machine Learning

## Possible queries



- Inference *← given some evidence observable*  
 $P(F=t \mid N=t)$   
*← infer state of some variable*
- Most probable explanation  
 $\max_{f,a,s,h} P(f,s,a,h \mid N=t)$
- Active data collection  
*What variable should I observe next?*  
 $H=? \quad S=?$

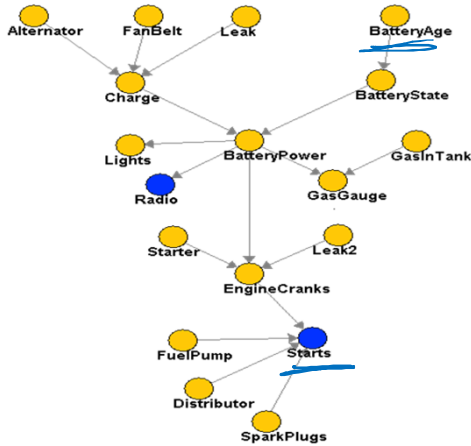
16

©2017 Emily Fox

CSF 446: Machine Learning



# CarStarts? Bayesian network



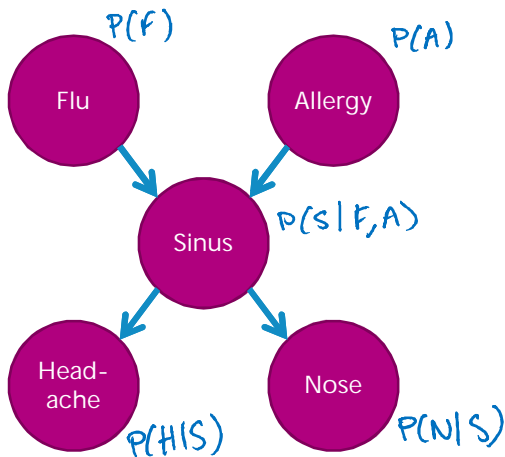
- 18 binary attributes  
*2<sup>18</sup> probabilities*
- Inference  
 -  $P(\text{BatteryAge} | \text{Starts}=f) = \frac{P(\text{BA}, S=f)}{P(S=f)} \propto P(\text{BA}, S=f)$   
 $= \sum P(a, f, l, c, \dots, S=f)$   
*(everything other than BatteryAge and Starts)*  
 (with *16* terms under the sum and *2<sup>16</sup>* next to it)
- 2<sup>16</sup> terms, why so fast?
- Not impressed?  
 - HailFinder BN – more than 3<sup>54</sup> = 58149737003040059690390169 terms

17

©2017 Emily Fox

CSF 446: Machine Learning

# Factored joint distribution – A preview



*2<sup>5</sup> = 32 terms... 31 params*

$$P(F, A, S, H, N)$$

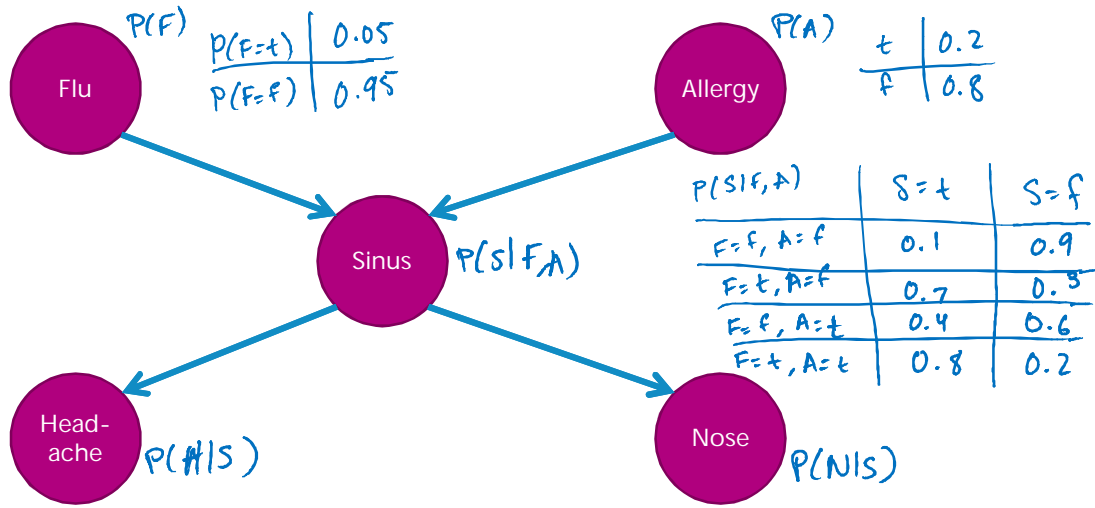
$$= P(F)P(A)P(S|F, A)P(H|S)P(N|S)$$

*(will see later why this fact. holds)*

©2017 Emily Fox

CSF 446: Machine Learning

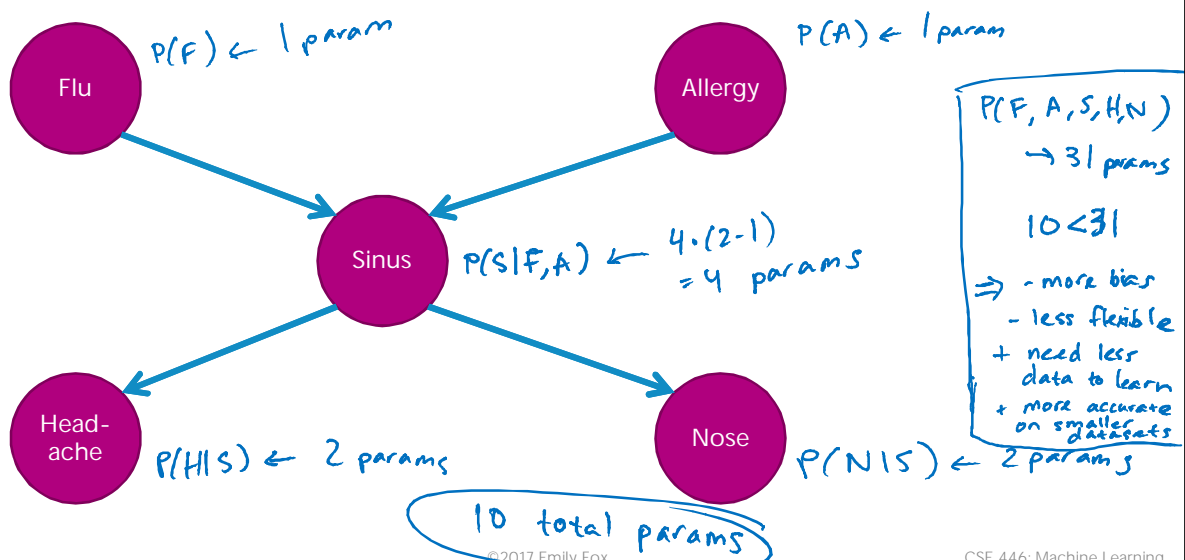
# What are these probabilities? Conditional probability tables (CPTs)



©2017 Emily Fox

CSF 446: Machine Learning

# Number of parameters



©2017 Emily Fox

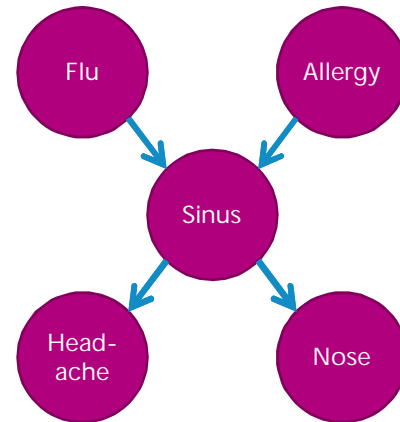
CSF 446: Machine Learning

## Factorization speeds up inference

Exploit distributivity:

$$\begin{aligned}
 P(F = x_F | N = t) &\propto \sum_{x_A, x_S, x_H} P(F = x_F, A = x_A, S = x_S, H = x_H, N = t) \\
 &= \sum_{x_A, x_S, x_H} P(F = x_F) P(A = x_A) P(S = x_S | F = x_F, A = x_A) P(H = x_H | S = x_S) P(N = t | S = x_S) \\
 &= P(F = x_F) \sum_{x_A} P(A = x_A) \sum_{x_S} P(S = x_S | F = x_F, A = x_A) P(N = t | S = x_S) \sum_{x_H} P(H = x_H | S = x_S)
 \end{aligned}$$

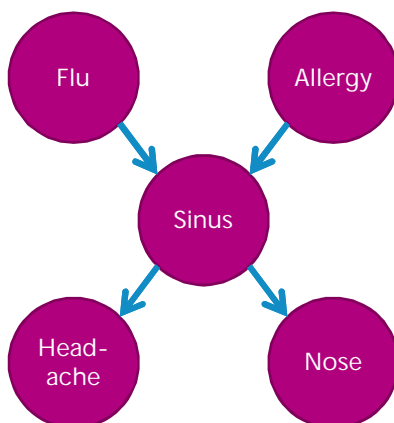
*variable elimination algorithm*



©2017 Emily Fox

CSF 446: Machine Learning

## Key: Independence assumptions



*FLN | S*

- Flu only "causes" Nose through Sinus
- If you tell  $N=t$ , this changes prob. of  $F$ , but if I first tell you  $S=t$ ,  $N$  doesn't affect prob. of  $F$

Knowing sinus separates variables from each other

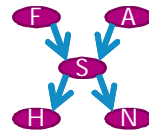
©2017 Emily Fox

CSF 446: Machine Learning

## Marginal and conditional independence

### (Marginal) Independence

- Flu and Allergy are (marginally) independent



$$F \perp A$$

$$P(A, F) = P(A)P(F)$$

$$\Updownarrow$$

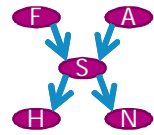
$$P(A|F) = P(A)$$

$P(F)$	Flu = t	0.2
	Flu = f	0.8

$P(A)$	Allergy = t	0.4
	Allergy = f	0.6

$P(A, F)$		Flu = t	Flu = f
	Allergy = t	$0.4 \times 0.2 = 0.08$	$0.4 \times 0.8$
	Allergy = f	$0.8 \times 0.2$	$0.8 \times 0.6$

## Conditional independence



- **Flu** and **Headache** are not (marginally) ind.  
 $P(H=t | F=t) \neq P(H=t)$
- **Flu** and **Headache** are independent given **Sinus** infection  
 $P(H=t | F=t, S=t) = P(H=t | S=t)$
- More generally:

$$X \perp Y | Z$$

$$P(X, Y | Z) = P(X | Z) P(Y | Z)$$

$$P(X | Y, Z) \stackrel{\updownarrow}{=} P(X | Z)$$

25

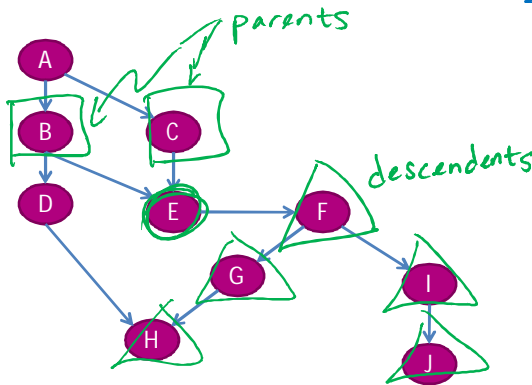
©2017 Emily Fox

CSF 446: Machine Learning

Conditional independence statements  
 encoded by Bayesian networks

# What is a Bayes net assuming?

**Local Markov Assumption:** A variable  $X$  is independent of its non-descendants given its parents



$$E \perp A \mid B, C$$

$$E \perp D \mid B, C$$

$$F \perp B \mid E$$

Allows you to read off some simple conditional independence relationships

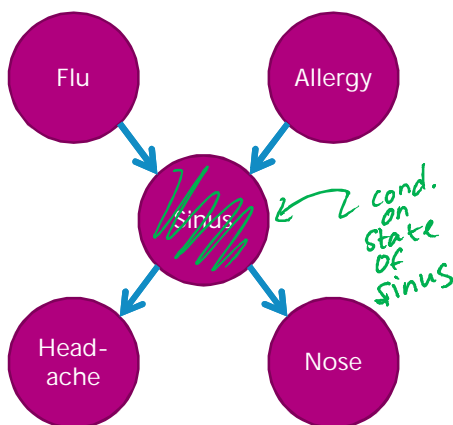
27

©2017 Emily Fox

CSF 446: Machine Learning

# Explaining away example

**Local Markov Assumption:**  
A variable  $X$  is independent of its non-descendants given its parents



$F \perp A$  (cond. ind. given  $\phi = \text{margin. ind.}$ )

$F \perp A \mid S$  ?? don't know

$P(F=t \mid A=t, S=t) \stackrel{?}{=} P(F=t \mid S=t)$

Suppose  $P(F=t \mid S=t)$  is high  
but  $P(F=t \mid S=t, A=t)$  is lower  
because  $A=t$  explains away  
sinus inflammation.

28

©2017 Emily Fox

CSF 446: Machine Learning

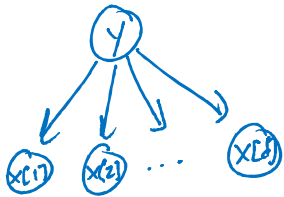
## Naïve Bayes revisited

**Local Markov Assumption:**  
A variable  $X$  is independent of its non-descendants given its parents

$$X^{[1]} \perp X^{[2]}, \dots, X^{[d]} \mid y$$

↑ true for any feature and set of remaining features

$$P(y, X^{[1]}, \dots, X^{[d]}) = P(y) \prod_{j=1}^d P(X^{[j]} \mid y) \quad \text{Naïve Bayes}$$



Local Markov assumption of this graph:

$$X^{[j]} \perp x^{[1]}, \dots, x^{[j-1]}, x^{[j+1]}, \dots, x^{[d]} \mid y$$

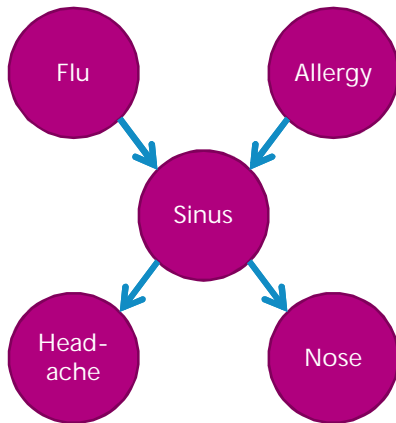
← Naïve Bayes!

©2017 Emily Fox

CSE 446: Machine Learning

## Factorization of the joint distribution

## Joint distribution



Why can we decompose?  
**Markov Assumption!**

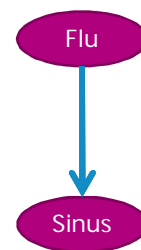
31

©2017 Emily Fox

CSF 446: Machine Learning

## The chain rule of probabilities

- $P(A,B) = P(A)P(B|A)$



- More generally:

- $P(\mathbf{X}[1], \dots, \mathbf{X}[d]) = P(\mathbf{X}[1]) P(\mathbf{X}[2]|\mathbf{X}[1]) \dots P(\mathbf{X}[d]|\mathbf{X}[1], \dots, \mathbf{X}[d-1])$

32

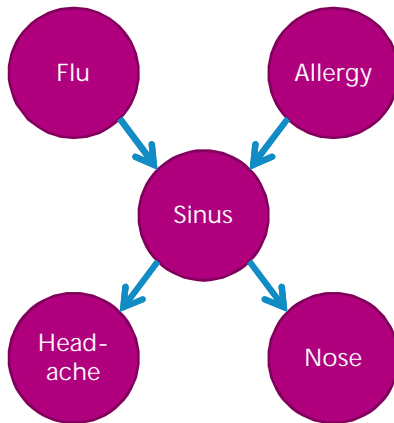
©2017 Emily Fox

CSF 446: Machine Learning



## Chain rule & joint distribution

**Local Markov Assumption:**  
A variable  $X$  is independent of its non-descendants given its parents

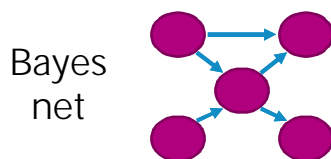


Order of expansion matters! Use topological order

©2017 Emily Fox

CSF 446: Machine Learning

## The Representation Theorem – Joint distribution to BN



Encodes  
independence  
assumptions

If cond. ind. in Bayes net  
are subset of cond. ind. in  $P$



Joint distribution factorizes:

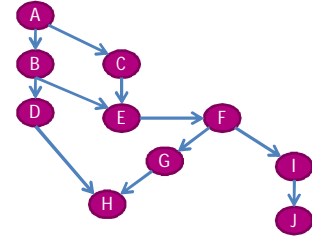
$$P(\mathbf{X}) = \prod_{j=1}^d P(\mathbf{X}[j] \mid \text{parents}(\mathbf{X}[j]))$$

34

©2017 Emily Fox

CSF 446: Machine Learning

## Bayesian networks recap



- Representation benefits
  - Compact representation for probability distributions
  - Exponential reduction in number of parameters
  - Lower variance parameter estimates from limited data
- Inference benefits
  - Efficient computation of  $P(X|e)$  (i.e., fast probabilistic inference)
  - Involves variable elimination algorithms
- Other important topics
  - Structure learning: What graph structure to use?
  - Understanding how evidence can be incorporated and how this changes conditional independence statements (d separation)

35

©2017 Emily Fox

CSF 446: Machine Learning

Hidden Markov models:  
A Bayesian network for time series

# Example: Motion Capture Segmentation



37

©2017 Emily Fox

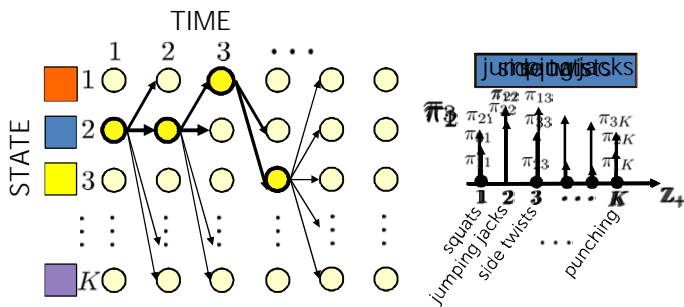
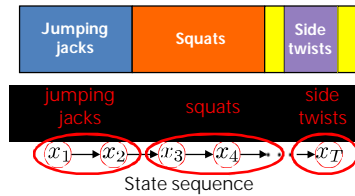
CSF 446: Machine Learning

# Hidden Markov model

Tutorial: Rabiner, Proc. IEEE 1989

Markov transition dynamics:

$$\Pr(x_t = \text{orange} \mid x_{t-1} = \text{blue}) = A_{\text{orange, blue}}$$



$$A = \begin{matrix} & \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ K \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ K \end{matrix} & \begin{matrix} \text{green} \\ \text{green} \\ \text{green} \\ \text{blue} \end{matrix} \end{matrix}$$

38

©2017 Emily Fox

CSF 446: Machine Learning

# Hidden Markov model

Tutorial: Rabiner, Proc. IEEE 1989

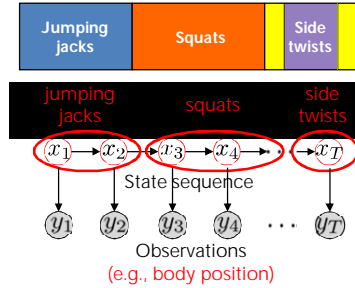
Markov transition dynamics:

$$\Pr(x_t = \blacksquare | x_{t-1} = \blacksquare) = A_{\blacksquare}$$

Conditionally independent emissions:

$$\Pr(y_t | x_t = \blacksquare) = N(\mu_{\blacksquare}, \Sigma_{\blacksquare})$$

Joint distribution factorization:



Latent Markov chain structure enables

- Efficient computation of marginals  $p(x_t | y_1, \dots, y_T)$  via **forward-backward** alg.
- Most-probable sequence via **Viterbi**
- Parameter learning via **Baum-Welch** (EM for HMMs)

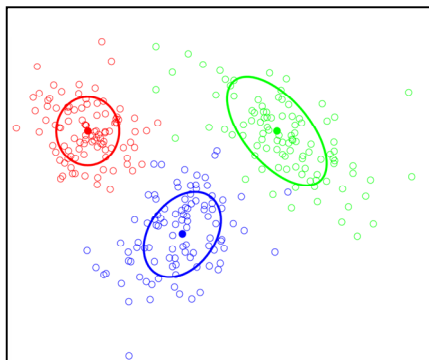
39

©2017 Emily Fox

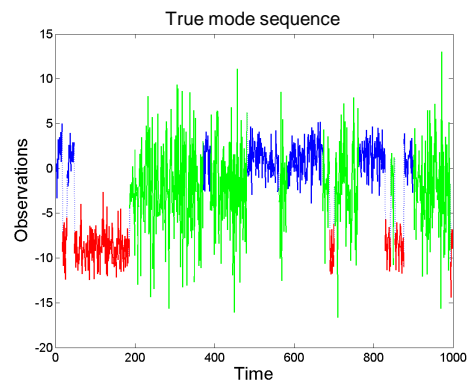
CSF 446: Machine Learning

# GMMs vs. HMMs

Gaussian mixture model



Hidden Markov model



40

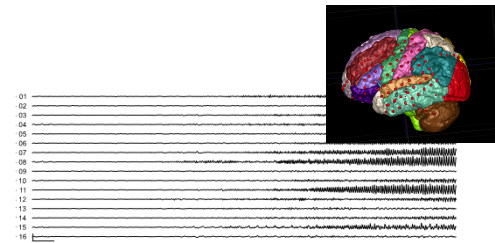
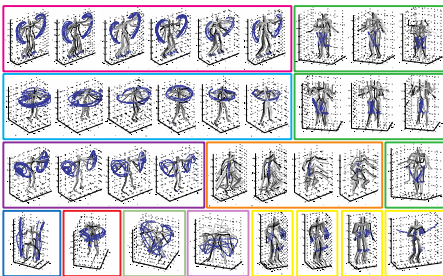
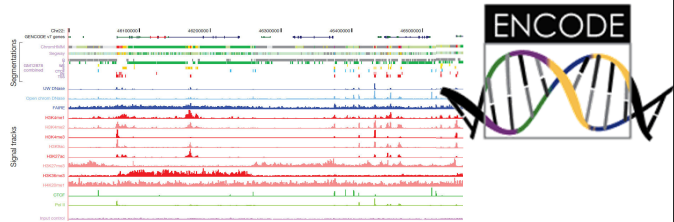
©2017 Emily Fox

CSF 446: Machine Learning

# HMM applications

## Example applications:

- Parsing EEG recordings
- Discovering behaviors in videos
- Speech segmentation
- Volatility regimes in financial time series
- Genomics
- ...



41

©2017 Emily Fox

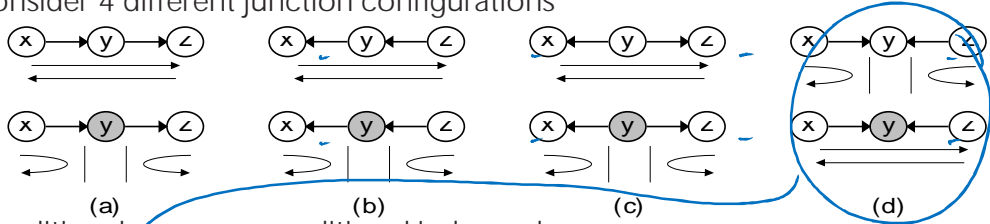
CSE 446: Machine Learning

Incorporating evidence:  
Bayes ball algorithm for analyzing  
conditional independencies

**OPTIONAL**

# Conditional independence in Bayes nets

- Consider 4 different junction configurations



- Conditional versus unconditional independence:

$P(x, y, z) = P(x)P(z)P(y|x, z) \stackrel{\text{int. over } y}{\Rightarrow} P(x, z) = P(x)P(z) \Rightarrow x \perp\!\!\!\perp z$   
 $P(x, z|y) \neq P(x, y, z) = P(x)P(z)P(y|x, z) \neq P(x|y)P(z|y) \leftarrow x \not\perp\!\!\!\perp z | y$   
 "explaining away":  $x = \text{earthquake}$ ,  $z = \text{burglar}$ ,  $y = \text{car alarm}$   
 If alarm ( $y=1$ ), an increase in earthquake  $p(x|y)$ , means  $p(z|y)$  lower  
ind. a priori

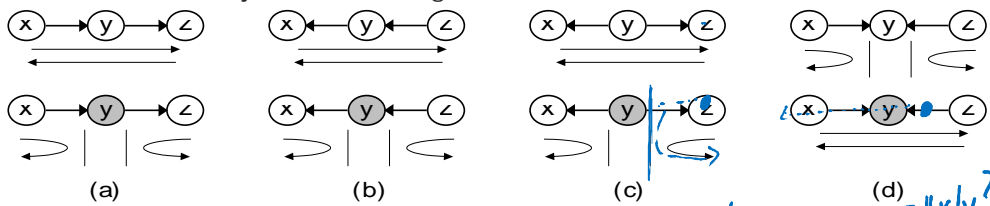
43

©2017 Emily Fox

CSF 446: Machine Learning

# Bayes ball algorithm

- Consider 4 different junction configurations



- Bayes ball algorithm:

Start ball at one end or other.  
 If ball passes to a node (straight arrows) then, not cond./marg. ind.  
 If ball bounces back (walls + curved arrows), the nodes are cond./marg. ind.

yes!  
 $z \perp\!\!\!\perp x | y$ ?  
 $\rightarrow$  No!

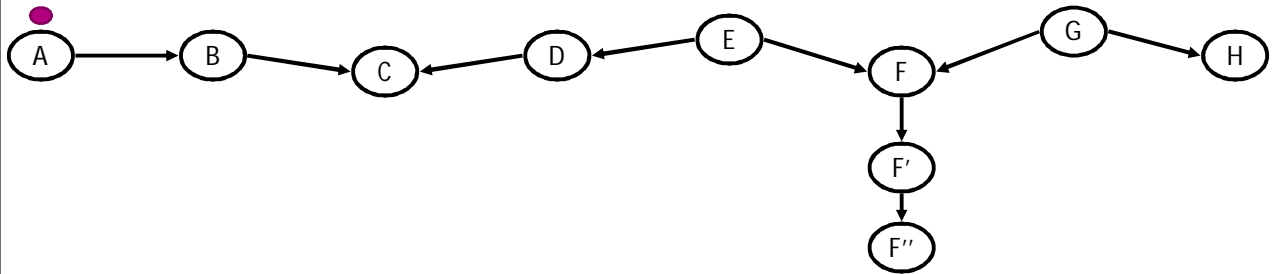
44

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



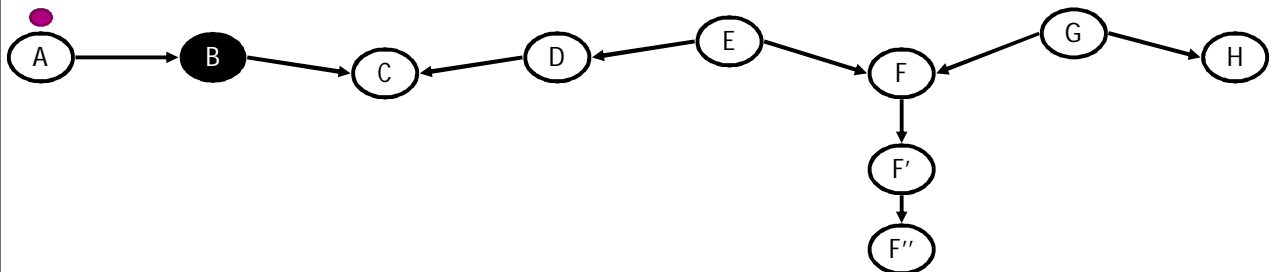
45

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



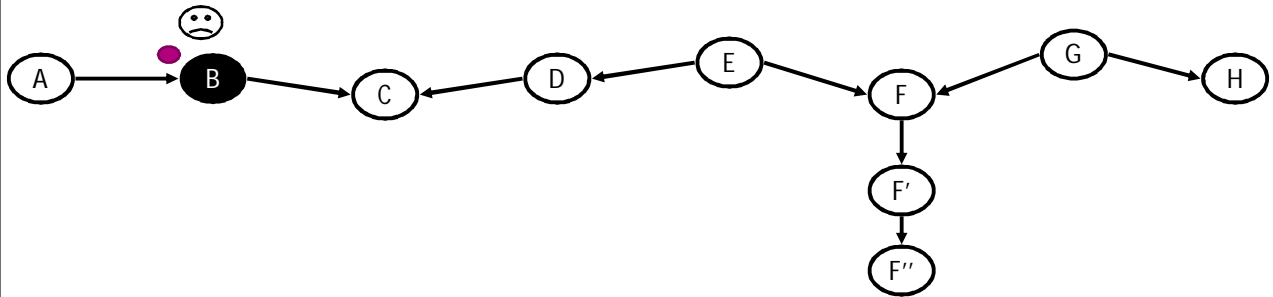
46

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



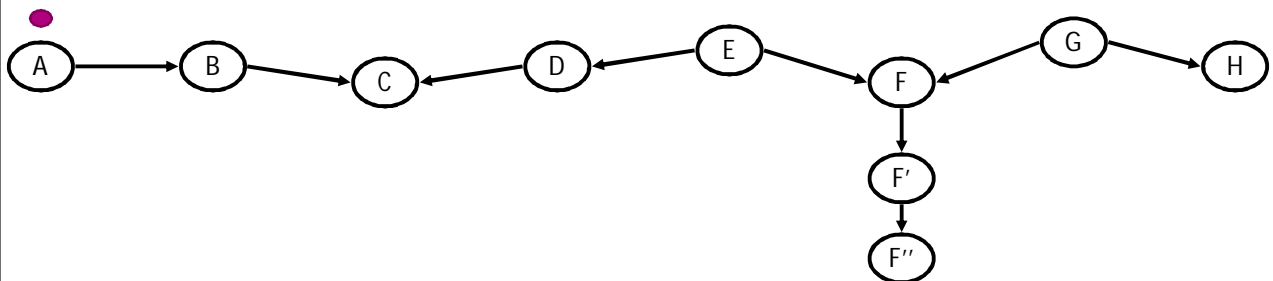
47

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



48

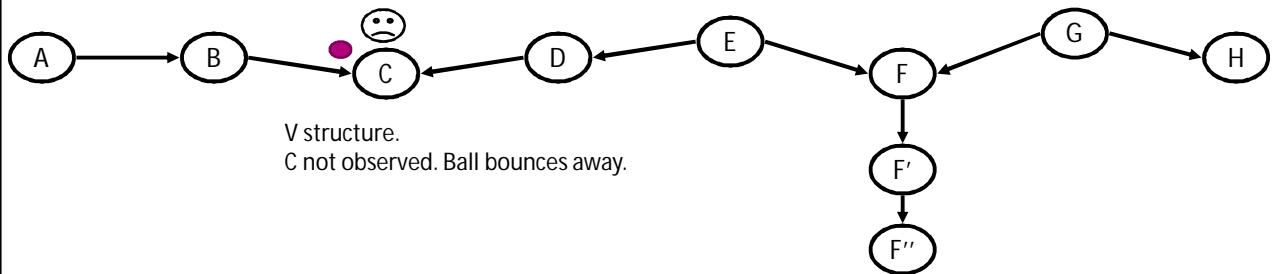
©2017 Emily Fox

CSF 446: Machine Learning



## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



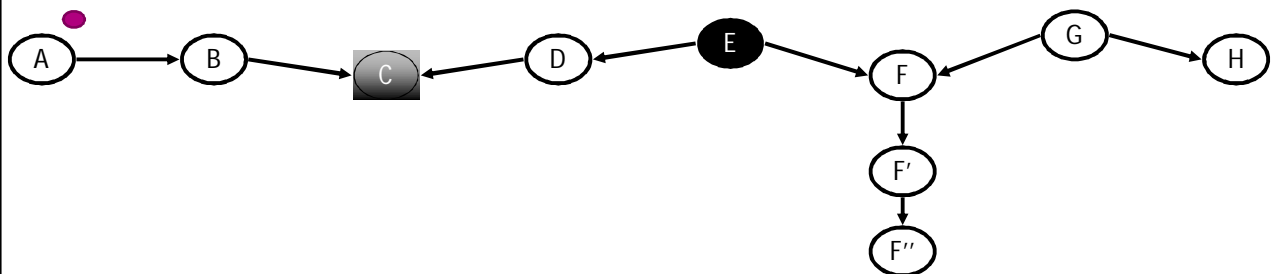
49

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



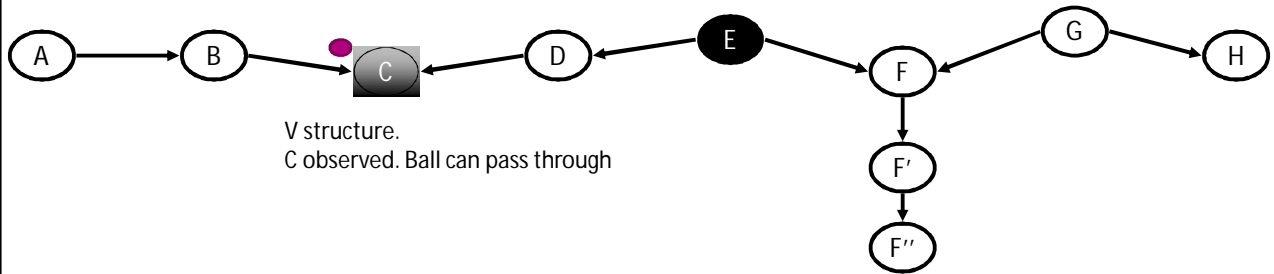
50

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



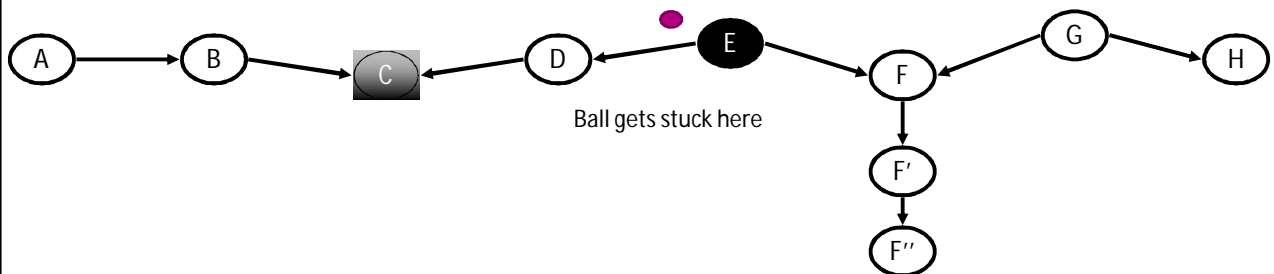
51

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



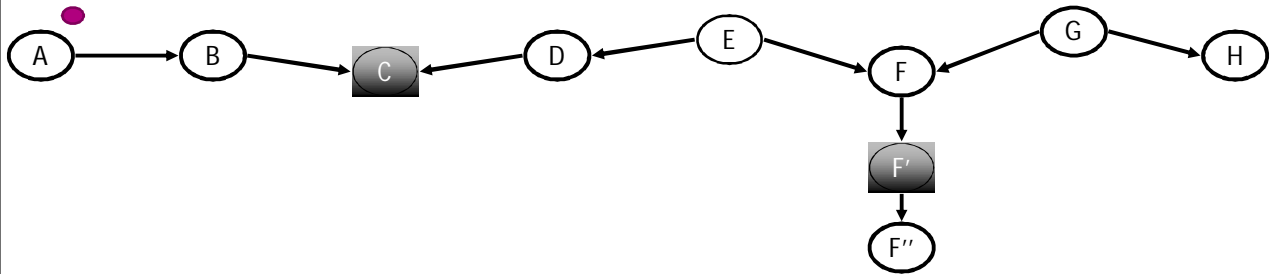
52

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



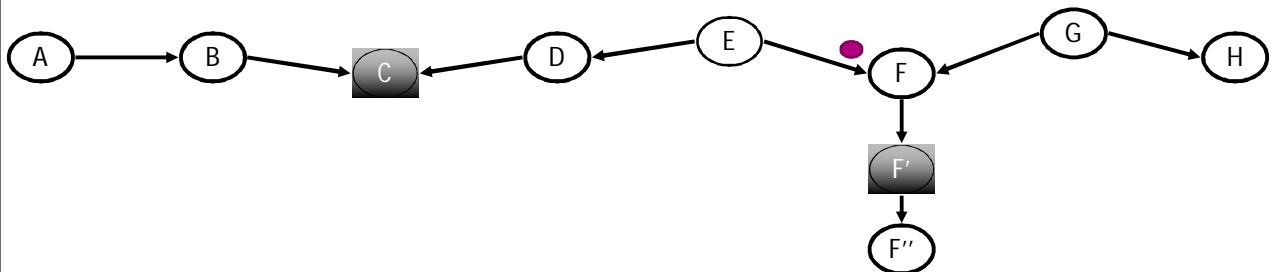
53

©2017 Emily Fox

CSF 446: Machine Learning

## Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H



V structure.  
Descendent of F observed.  
Ball can pass through

54

©2017 Emily Fox

CSF 446: Machine Learning

# Bayes ball example

A path from A to H is Active if the Bayes ball can get from A to H

