

# Bayes Optimal Classifier & Naïve Bayes

CSE 446: Machine Learning  
Emily Fox  
University of Washington  
March 3, 2017

©2017 Emily Fox

## Classification

Learn:  $f: \mathbf{X} \mapsto Y$    
*← GPA, 446 grade, ...*  
*← hired / not hired*

- $\mathbf{X}$  – features
- $Y$  – target classes

Suppose you know  $P(Y|\mathbf{X})$  exactly, how should you classify?

- Bayes optimal classifier:

$$\hat{y} = \arg \max_y P(Y=y | X=x)$$

$X = x = \{ \text{GPA}=3.8, \text{446 grade}=3.9 \}$

In logistic regression, model for  $P(Y|X=x) = \frac{1}{1+e^{-w \cdot x}}$   
 Model  $P(Y|X)$  directly.. "discriminative model"

## Recall: Bayes rule

$$P(Y | \mathbf{X}) = \frac{P(\mathbf{X} | Y)P(Y)}{P(\mathbf{X})}$$

"likelihood"
"prior"
normalizer

Which is shorthand for:

$$(\forall i, \mathbf{j}) \quad P(Y = i | \mathbf{X} = \mathbf{j}) = \frac{P(\mathbf{X} = \mathbf{j} | Y = i)P(Y = i)}{P(\mathbf{X} = \mathbf{j})}$$

3

©2017 Emily Fox

CSE 446: Machine Learning

## How hard is it to learn the optimal classifier?

How hard is it to learn the output?

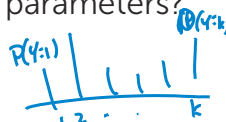
|          | $x[1]$ | $x[2]$ | ...    | $x[d]$ | $Y$   |          |          |
|----------|--------|--------|--------|--------|-------|----------|----------|
|          | Sky    | Temp   | Humid  | Wind   | Water | Forecast | EnjoySpt |
| • Data = | Sunny  | Warm   | Normal | Strong | Warm  | Same     | Yes      |
|          | Sunny  | Warm   | High   | Strong | Warm  | Same     | Yes      |
|          | Rainy  | Cold   | High   | Strong | Warm  | Change   | No       |
|          | Sunny  | Warm   | High   | Strong | Cold  | Change   | Yes      |

- How do we represent these? How many parameters?

- Prior,  $P(Y)$ :

- Suppose  $Y$  is composed of  $k$  classes

$k-1$



$$\sum P(Y=y) = 1$$

- Likelihood,  $P(\mathbf{X}|Y)$ :

- Suppose  $\mathbf{X}$  is composed of  $d$  binary features

$P(\mathbf{X}=\mathbf{x} | Y=y) \leftarrow$  for each class ( $Y=y$ ), dist over feature values

$k(2^d - 1)$   $\leftarrow$  a lot of params! need a lot of data,

- Complex model ! High variance with limited data!!!

4

©2017 Emily Fox

CSE 446: Machine Learning

## Conditional Independence

X is **conditionally independent** of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z

$$(\forall i, j, k) \quad P(X = i \mid Y = j, Z = k) = P(X = i \mid Z = k)$$

e.g.,

$$P(\text{Thunder} \mid \text{Rain}, \text{Lightening}) = P(\text{Thunder} \mid \text{Lightening})$$

Equivalent to:

$$P(X, Y \mid Z) = P(X \mid Z)P(Y \mid Z)$$

5

©2017 Emily Fox

CSE 446: Machine Learning

## What if features are independent?

- Predict **Lightening**
- From two conditionally independent features
  - **Thunder**
  - **Rain**

6

©2017 Emily Fox

CSE 446: Machine Learning

## The Naïve Bayes assumption

- Naïve Bayes assumption:
  - Features are independent given class:

$$\begin{aligned} P(\mathbf{X}[1], \mathbf{X}[2] \mid Y) &= P(\mathbf{X}[1] \mid \mathbf{X}[2], Y)P(\mathbf{X}[2] \mid Y) \\ &= P(\mathbf{X}[1] \mid Y)P(\mathbf{X}[2] \mid Y) \end{aligned}$$

- More generally:

$$P(\mathbf{X}[1], \dots, \mathbf{X}[d] \mid Y) = \prod_j P(\mathbf{X}[j] \mid Y)$$

- How many parameters now?
  - Suppose  $\mathbf{X}$  is composed of  $d$  binary features

7

©2017 Emily Fox

CSE 446: Machine Learning

## The Naïve Bayes classifier

- Given:
  - Prior  $P(Y)$
  - $d$  conditionally independent features  $\mathbf{X}[j]$  given the class  $Y$
  - For each  $\mathbf{X}[j]$ , we have likelihood  $P(\mathbf{X}[j] \mid Y)$

- Decision rule:

$$\begin{aligned} \hat{y} = f_{NB}(\mathbf{x}) &= \arg \max_y P(y)P(\mathbf{x}[1], \dots, \mathbf{x}[d] \mid y) \\ &= \arg \max_y P(y) \prod_j P(\mathbf{x}[j] \mid y) \end{aligned}$$

- If assumption holds, NB is optimal classifier!

8

©2017 Emily Fox

CSE 446: Machine Learning



## MLE for the parameters of NB

- Given dataset
  - $\text{Count}(A=a, B=b) == \# \text{ examples where } A=a \text{ and } B=b$
- MLE for NB, simply:
  - Prior:  $P(Y=y) =$
  - Likelihood:  $P(\mathbf{X}[j]=\mathbf{x}[j] \mid Y=y) =$

9

©2017 Emily Fox

CSE 446: Machine Learning

## Subtleties of NB classifier 1 – Violating the NB assumption

- Usually, features are not conditionally independent:

$$P(\mathbf{X}[1], \dots, \mathbf{X}[d] \mid Y) \neq \prod_j P(\mathbf{X}[j] \mid Y)$$

- Actual probabilities  $P(Y|\mathbf{X})$  often biased towards 0 or 1
- Nonetheless, NB is one of the most used classifier out there
  - NB often performs well, even when assumption is violated
  - [Domingos & Elkan '96] discuss some conditions for good performance

10

©2017 Emily Fox

CSE 446: Machine Learning

## Subtleties of NB classifier 2 – Insufficient training data

- What if you never see a training instance where  $X[1]=a$  when  $Y=b$ ?
  - e.g.,  $Y=\{\text{SpamEmail}\}$ ,  $X[1]=\{\text{'Viagra'}\}$
  - $P(X[1]=a \mid Y=b) = 0$
- Thus, no matter what the values  $X[2], \dots, X[d]$  take:
  - $P(Y=b \mid X[1]=a, X[2], \dots, X[d]) = 0$
- “Solution”: **smoothing**
  - Add “fake” counts, usually uniformly distributed
  - Equivalent to **Bayesian learning**

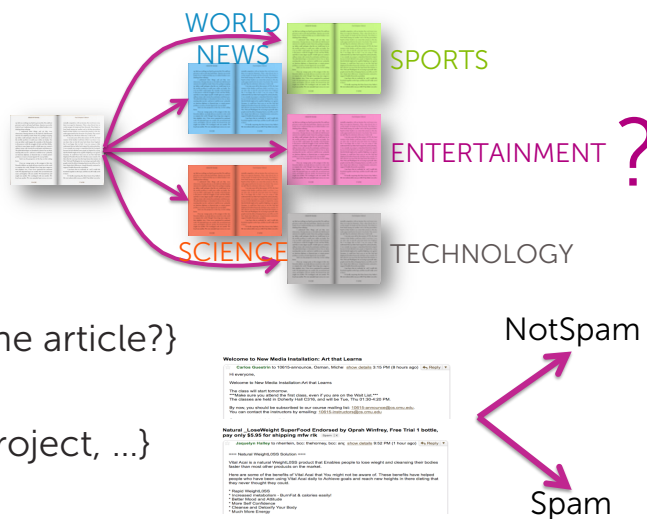
11

©2017 Emily Fox

CSE 446: Machine Learning

## Text classification

- Classify e-mails
  - $Y = \{\text{Spam}, \text{NotSpam}\}$
- Classify news articles
  - $Y = \{\text{what is the topic of the article?}\}$
- Classify webpages
  - $Y = \{\text{Student, professor, project, ...}\}$
- What about the features  $X$ ?
  - The text!



12

©2017 Emily Fox

CSE 446: Machine Learning

## Features $\mathbf{X}$ are entire document – $\mathbf{X}[j]$ for $j$ th word in article

Article from rec.sport.hockey

---

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e  
 From: xxx@yyy.zzz.edu (John Doe)  
 Subject: Re: This year's biggest and worst (opinion)  
 Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

13

CSE 446: Machine Learning

## NB for text classification

- $P(\mathbf{X}|Y)$  is huge!!!
  - Article at least 1000 words,  $\mathbf{X}=\{\mathbf{X}[1], \dots, \mathbf{X}[1000]\}$
  - $\mathbf{X}[j]$  represents  $j$ th word in document
    - i.e., the domain of  $\mathbf{X}[j]$  is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.
- NB assumption helps a lot!!!
  - $P(\mathbf{X}[j]=\mathbf{x}[j]|Y=y)$  is the probability of observing word  $\mathbf{x}[j]$  in a document on topic  $y$

$$f_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{j=1}^{LengthDoc} P(\mathbf{x}[j] | y)$$

14

©2017 Emily Fox

CSE 446: Machine Learning

## Bag of words model

- Typical additional assumption: **Position in document doesn't matter**

$$P(\mathbf{X}[j]=\mathbf{x}[j] \mid Y=y) = P(\mathbf{X}[k]=\mathbf{x}[j] \mid Y=y)$$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{j=1}^{LengthDoc} P(\mathbf{x}[j] \mid y)$$

When the lecture is over, remember to wake up the person sitting next to you in the lecture room.

15

©2017 Emily Fox

CSE 446: Machine Learning

## Bag of words model

- Typical additional assumption: **Position in document doesn't matter**

$$P(\mathbf{X}[j]=\mathbf{x}[j] \mid Y=y) = P(\mathbf{X}[k]=\mathbf{x}[j] \mid Y=y)$$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{j=1}^{LengthDoc} P(\mathbf{x}[j] \mid y)$$

in is lecture lecture next over person remember room  
sitting the the the to to up wake when you

16

©2017 Emily Fox

CSE 446: Machine Learning

# Bag-of-words representation

Modeling the Complex Dynamics and Changing  
Correlations of Epileptic Events

Drausin F. Wulsin\*, Emily B. Fox<sup>a</sup>, Brian Litt<sup>a,b</sup>

<sup>a</sup>Department of Bioengineering, University of Pennsylvania, Philadelphia, PA  
<sup>b</sup>Department of Neurology, University of Pennsylvania, Philadelphia, PA  
<sup>c</sup>Department of Statistics, University of Washington, Seattle, WA

## Abstract

Patients with epilepsy can manifest short, sub-clinical epileptic “bursts” in addition to full-blown clinical seizures. We believe the relationship between these two classes of events—something not previously studied quantitatively—could yield important insights into the nature and intrinsic dynamics of seizures. A goal of our work is to parse these complex epileptic events into distinct dynamic regimes. A challenge posed by the intracranial EEG (IEEG) data we study is the fact that the number and placement of electrodes can vary between patients. We develop a Bayesian nonparametric Markov switching process that allows for (i) shared dynamic regimes between a variable number of channels, (ii) asynchronous regime-switching, and (iii) an unknown dictionary of dynamic regimes. We encode a sparse and changing set of dependencies between the channels using a Markov-switching Gaussian graphical model for the innovations process driving the channel dynamics and demonstrate the importance of this model in parsing and out-of-sample predictions of IEEG data. We show that our model produces intuitive state assignments that can help automate clinical analysis of seizures and enable the comparison of sub-clinical bursts and full clinical seizures.

**Keywords:** Bayesian nonparametric, EEG, factorial hidden Markov model, graphical model, time series

## 1. Introduction

Despite over three decades of research, we still have very little idea of what defines a seizure. This ignorance stems both from the complexity of epilepsy as a disease and a paucity of quantitative tools that are flexible



17

©2017 Emily Fox

CSE 446: Machine Learning

# Bag-of-words representation

Modeling the Complex Dynamics and Changing  
Correlations of Epileptic Events

Drausin F. Wulsin\*, Emily B. Fox<sup>a</sup>, Brian Litt<sup>a,b</sup>

<sup>a</sup>Department of Bioengineering, University of Pennsylvania, Philadelphia, PA  
<sup>b</sup>Department of Neurology, University of Pennsylvania, Philadelphia, PA  
<sup>c</sup>Department of Statistics, University of Washington, Seattle, WA

## Abstract

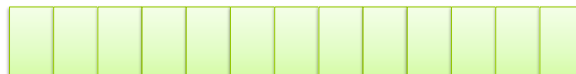
Patients with epilepsy can manifest short, sub-clinical epileptic “bursts” in addition to full-blown clinical seizures. We believe the relationship between these two classes of events—something not previously studied quantitatively—could yield important insights into the nature and intrinsic dynamics of seizures. A goal of our work is to parse these complex epileptic events into distinct dynamic regimes. A challenge posed by the intracranial EEG (IEEG) data we study is the fact that the number and placement of electrodes can vary between patients. We develop a Bayesian nonparametric Markov switching process that allows for (i) shared dynamic regimes between a variable number of channels, (ii) asynchronous regime-switching, and (iii) an unknown dictionary of dynamic regimes. We encode a sparse and changing set of dependencies between the channels using a Markov-switching Gaussian graphical model for the innovations process driving the channel dynamics and demonstrate the importance of this model in parsing and out-of-sample predictions of IEEG data. We show that our model produces intuitive state assignments that can help automate clinical analysis of seizures and enable the comparison of sub-clinical bursts and full clinical seizures.

**Keywords:** Bayesian nonparametric, EEG, factorial hidden Markov model, graphical model, time series

## 1. Introduction

Despite over three decades of research, we still have very little idea of what defines a seizure. This ignorance stems both from the complexity of epilepsy as a disease and a paucity of quantitative tools that are flexible

{modeling, complex, epilepsy,  
modeling, Bayesian, clinical,  
epilepsy, EEG, data, dynamic...}



18

©2017 Emily Fox

CSE 446: Machine Learning

## NB with bag of words for text classification

- Learning phase:
  - Prior  $P(Y)$ 
    - Count how many documents you have from each topic (+ prior)
  - $P(\mathbf{X}[j]|Y)$ 
    - For each topic, count how many times you saw word in documents of this topic (+ prior)
- Test phase:
  - For each document
    - Use naïve Bayes decision rule

$$f_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{j=1}^{LengthDoc} P(\mathbf{x}[j] | y)$$

19

©2017 Emily Fox

CSE 446: Machine Learning

## Twenty News Groups results

Given 1000 training documents from each group  
Learn to classify new documents according to  
which newsgroup it came from

|                          |                    |
|--------------------------|--------------------|
| comp.graphics            | misc.forsale       |
| comp.os.ms-windows.misc  | rec.autos          |
| comp.sys.ibm.pc.hardware | rec.motorcycles    |
| comp.sys.mac.hardware    | rec.sport.baseball |
| comp.windows.x           | rec.sport.hockey   |
| alt.atheism              | sci.space          |
| soc.religion.christian   | sci.crypt          |
| talk.religion.misc       | sci.electronics    |
| talk.politics.mideast    | sci.med            |
| talk.politics.misc       |                    |
| talk.politics.guns       |                    |

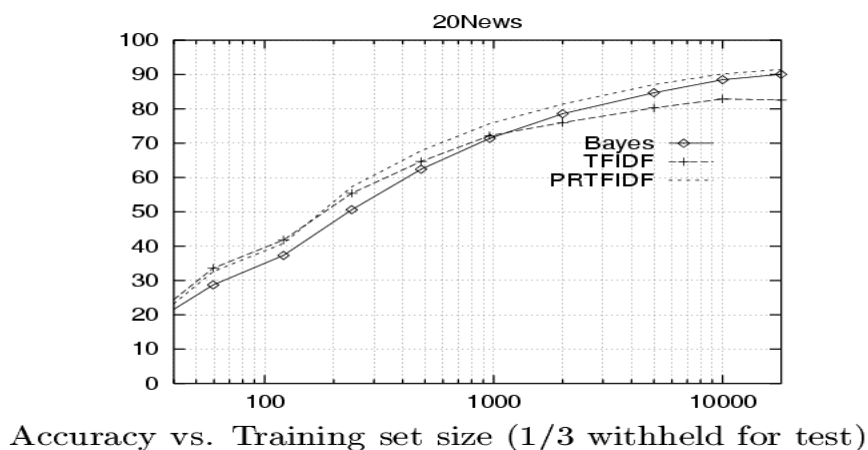
Naive Bayes: 89% classification accuracy

20

©2017 Emily Fox

CSE 446: Machine Learning

## Learning curve for Twenty News Groups



21

©2017 Emily Fox

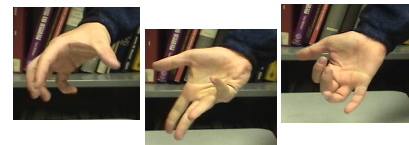
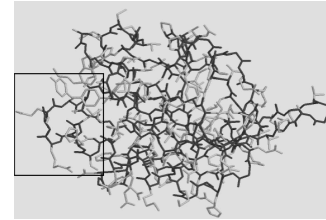
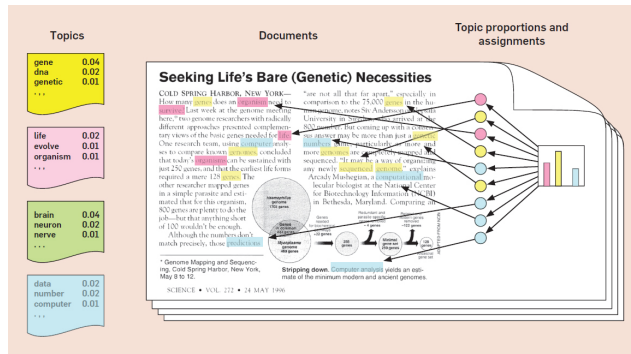
CSE 446: Machine Learning

## Bayesian Networks– Representation

CSE 446: Machine Learning  
Emily Fox  
University of Washington  
March 3, 2017

©2017 Emily Fox

# Learning from structured data

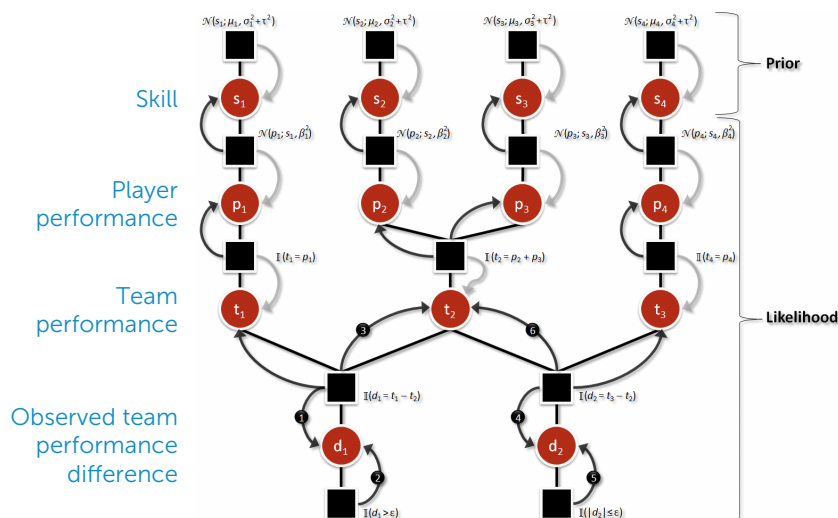


23

©2017 Emily Fox

CSE 446: Machine Learning

## TrueSkill: A Bayesian Skill Rating System



Herbrich et al., 2007

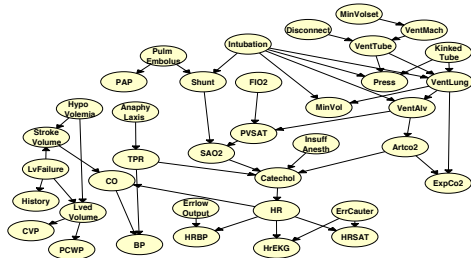
24

©2017 Emily Fox

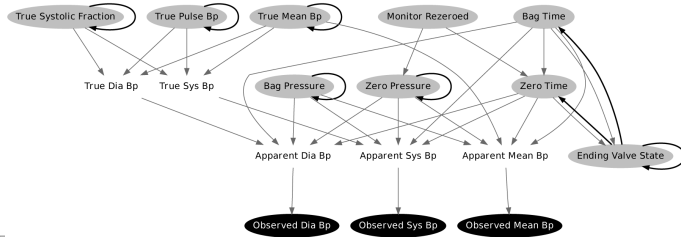
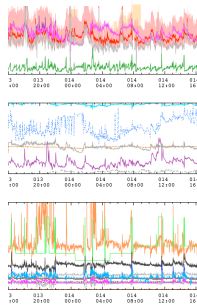
CSE 446: Machine Learning



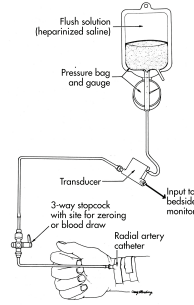
# ICU Monitoring



Beinlich et al., 1989



Aleks, Russell, et al., 2008



25

©2017 Emily Fox

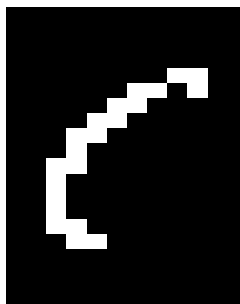
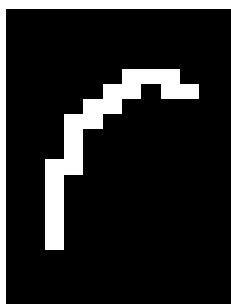
CSE 446: Machine Learning

Digging in:  
Learning with and without context/structure

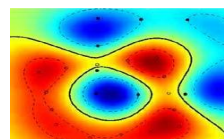
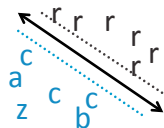
©2017 Emily Fox

CSE 446: Machine Learning

## Without context: Handwriting recognition



Character recognition,  
e.g., kernel SVMs

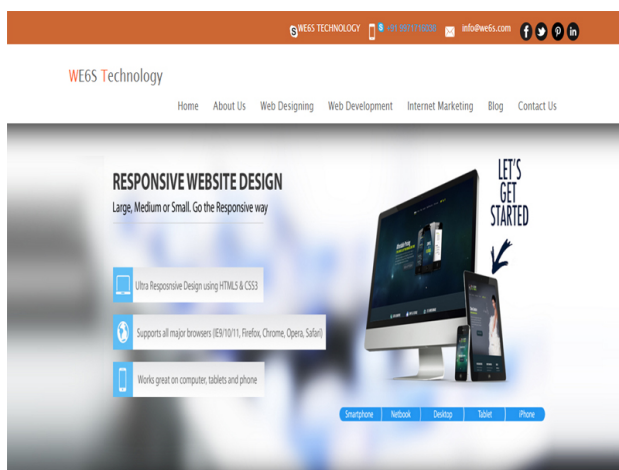


27

©2017 Emily Fox

CSE 446: Machine Learning

## Without context: Webpage classification



Company website

University website

Personal website

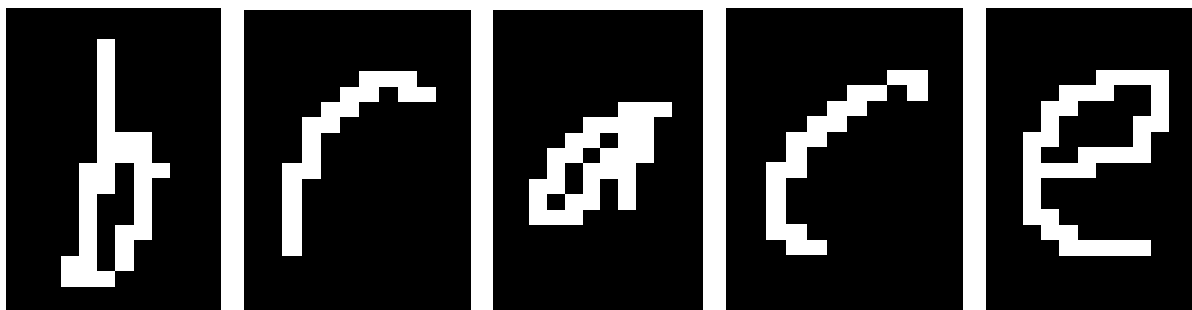
...

28

©2017 Emily Fox

CSE 446: Machine Learning

## With context: Handwriting recognition



©2017 Emily Fox

CSE 446: Machine Learning

## With context: Webpage classification



30

©2017 Emily Fox

CSE 446: Machine Learning

## Modeling structured relationships via Bayesian networks

©2017 Emily Fox

CSE 446: Machine Learning

### 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

32

©2017 Emily Fox

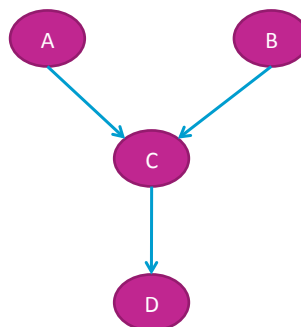
CSE 446: Machine Learning

## Bayesian network representation

Compact representation of a probability distribution.

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

Directed Acyclic Graph

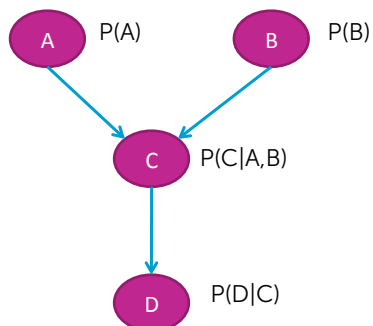


33

©2017 Emily Fox

CSE 446: Machine Learning

## Bayesian network probability factorization



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

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

implies the factorization:

$$P(X) = \prod_{i=1}^{|X|} P(X_i | \text{parents}(X_i))$$

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

34

©2017 Emily Fox

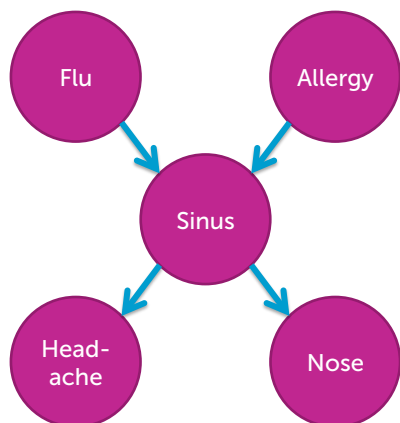
CSE 446: Machine Learning

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?

## Possible queries



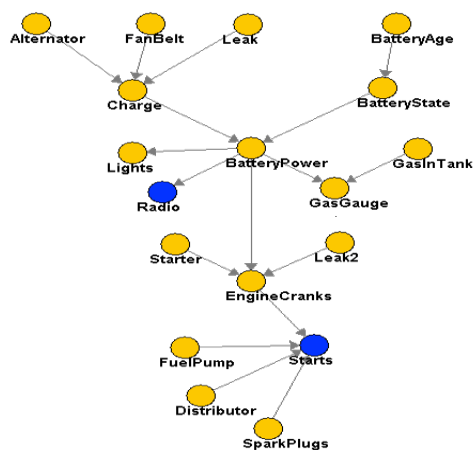
- Inference
- Most probable explanation
- Active data collection

37

©2017 Emily Fox

CSE 446: Machine Learning

## CarStarts? Bayesian network



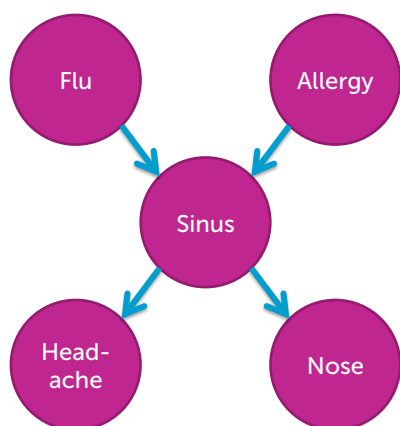
- 18 binary attributes
- Inference
  - $P(\text{BatteryAge} | \text{Starts}=f)$
- $2^{16}$  terms, why so fast?
- Not impressed?
  - HailFinder BN – more than  $3^{54} = 58149737003040059690390169$  terms

38

©2017 Emily Fox

CSE 446: Machine Learning

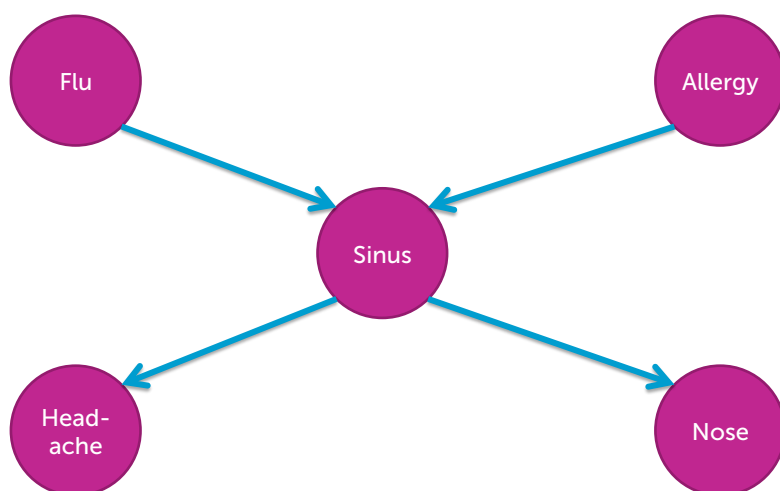
## Factored joint distribution – A preview



©2017 Emily Fox

CSE 446: Machine Learning

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

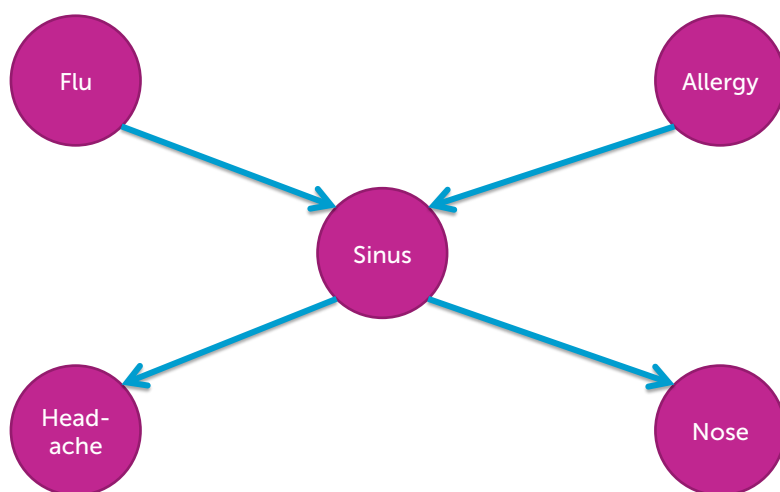


©2017 Emily Fox

CSE 446: Machine Learning



## Number of parameters



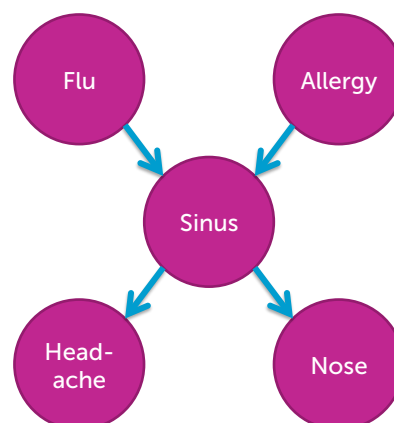
©2017 Emily Fox

CSE 446: Machine Learning

## Factorization speeds up inference

Exploit distributivity:

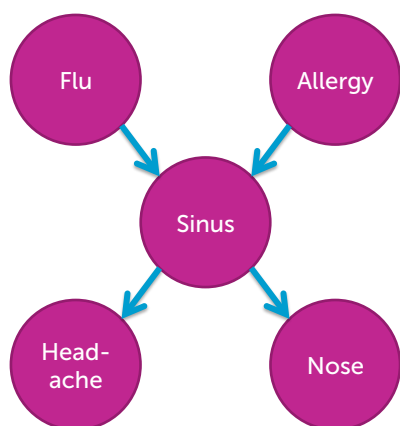
$$\begin{aligned}
 P(F = x_F | N = t) &\propto \sum_{x_A, x_S, x_H} P(F = x_F, A = x_A, S = x_S, H = x_H, N = t) \\
 &= \sum_{x_A, x_S, x_H} P(F = x_F) P(A = x_A) P(S = x_S | F = x_F, A = x_A) P(H = x_H | S = x_S) P(N = t | S = x_S) \\
 &= P(F = x_F) \sum_{x_A} P(A = x_A) \sum_{x_S} P(S = x_S | F = x_F,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



©2017 Emily Fox

CSE 446: Machine Learning

## Key: Independence assumptions



Knowing sinus **separates variables** from each other

©2017 Emily Fox

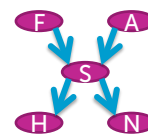
CSE 446: Machine Learning

## Marginal and conditional independence

©2017 Emily Fox

CSE 446: Machine Learning

## (Marginal) Independence



- Flu and Allergy are (marginally) independent

|         |  |
|---------|--|
| Flu = t |  |
| Flu = f |  |

|             |  |
|-------------|--|
| Allergy = t |  |
| Allergy = f |  |

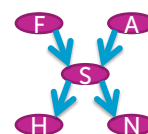
|             | Flu = t | Flu = f |
|-------------|---------|---------|
| Allergy = t |         |         |
| Allergy = f |         |         |

45

©2017 Emily Fox

CSE 446: Machine Learning

## Conditional independence



- Flu and Headache are not (marginally) ind.
- Flu and Headache are independent given Sinus infection
- More generally:

46

©2017 Emily Fox

CSE 446: Machine Learning

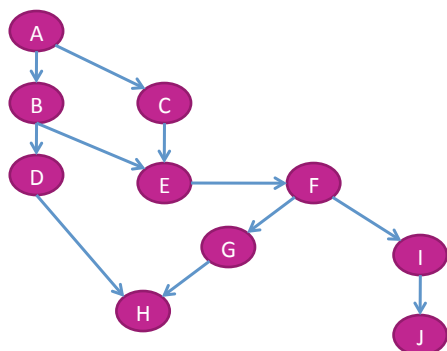
## Conditional independence statements encoded by Bayesian networks

©2017 Emily Fox

CSE 446: Machine Learning

## 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

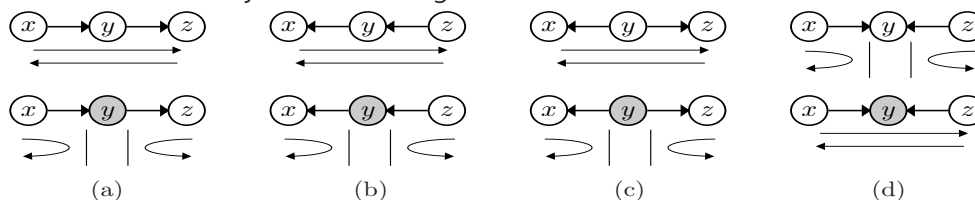
48

©2017 Emily Fox

CSE 446: Machine Learning

## Conditional independence in Bayes nets

- Consider 4 different junction configurations



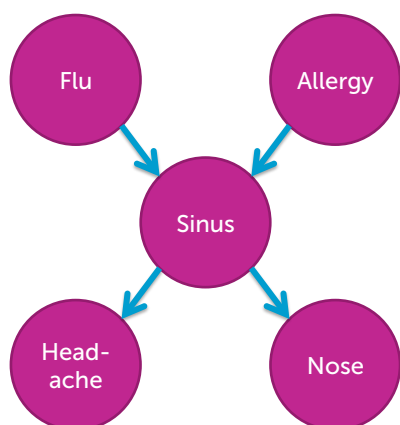
- Conditional versus unconditional independence:

49

©2017 Emily Fox

CSE 446: Machine Learning

## Explaining away example



### Local Markov Assumption:

A variable  $X$  is independent of its non-descendants given its parents

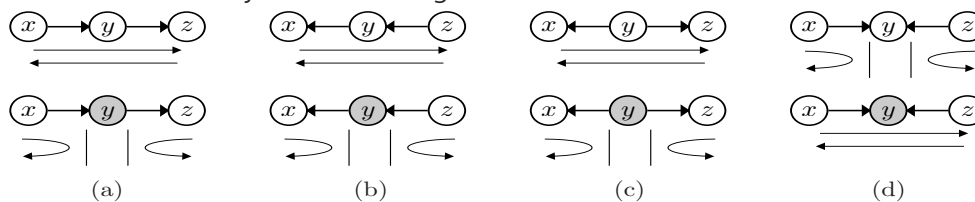
50

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball algorithm

- Consider 4 different junction configurations



- Bayes ball algorithm:

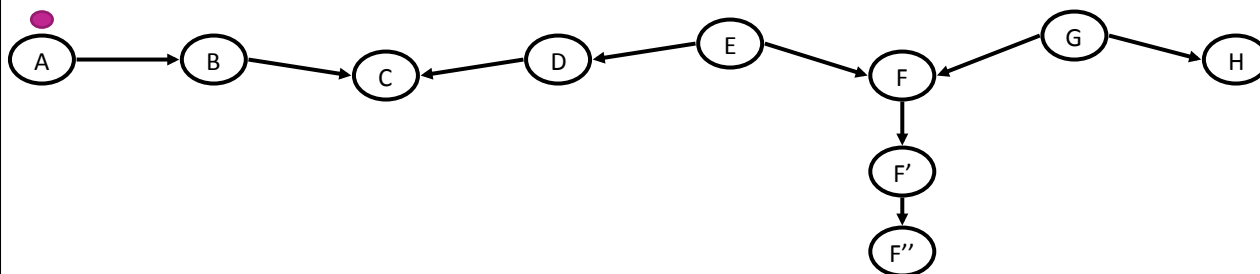
51

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



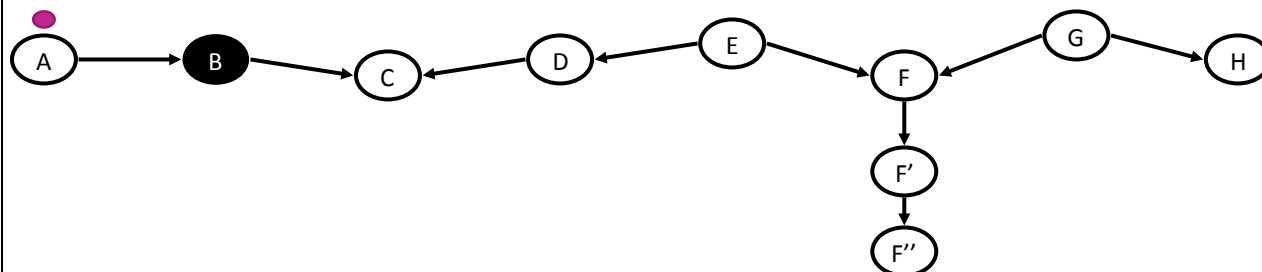
52

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



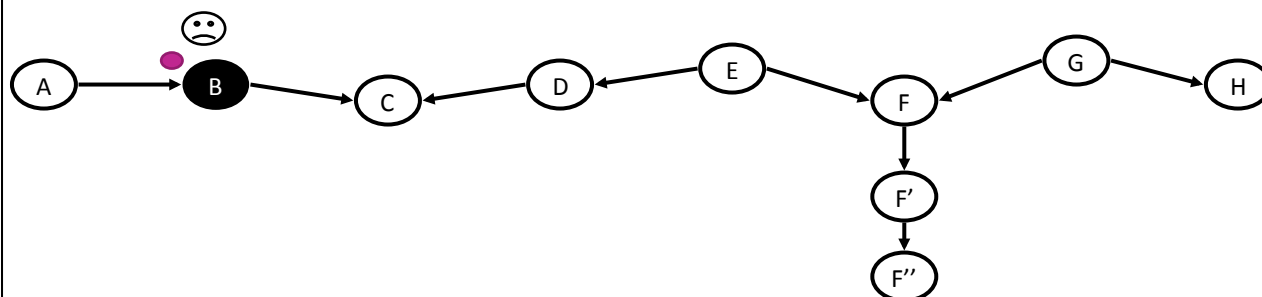
53

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



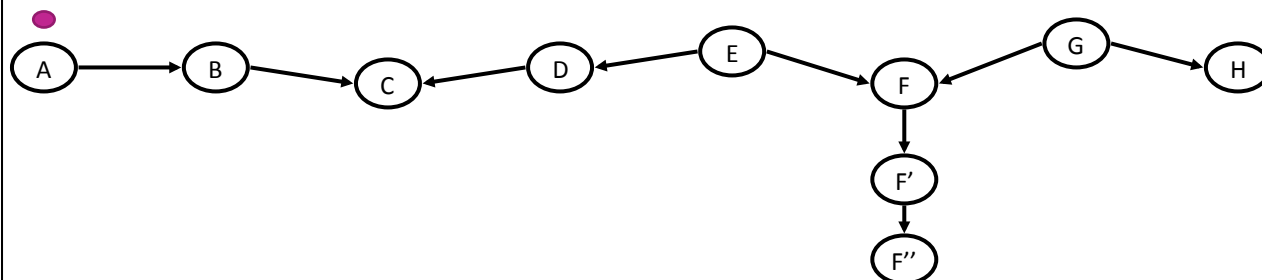
54

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



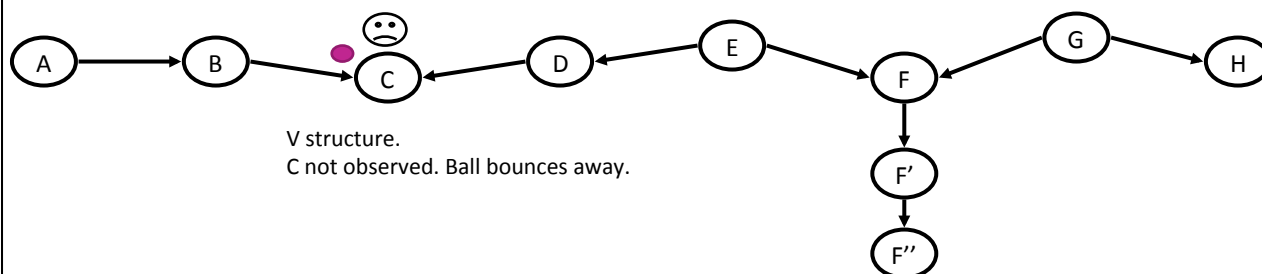
55

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



56

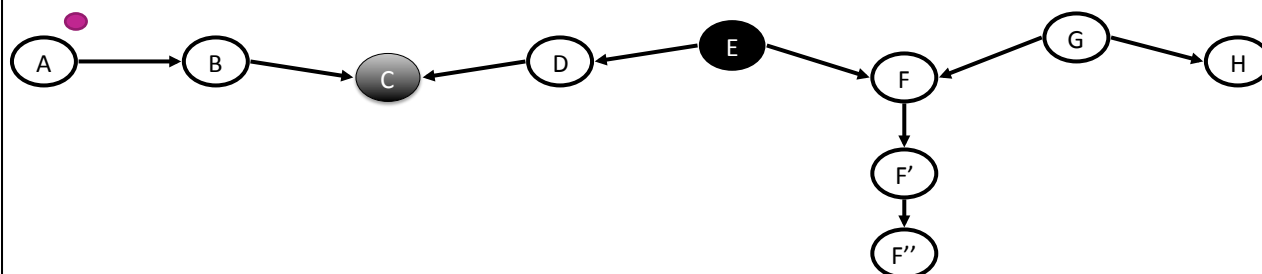
©2017 Emily Fox

CSE 446: Machine Learning



## Bayes ball example

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



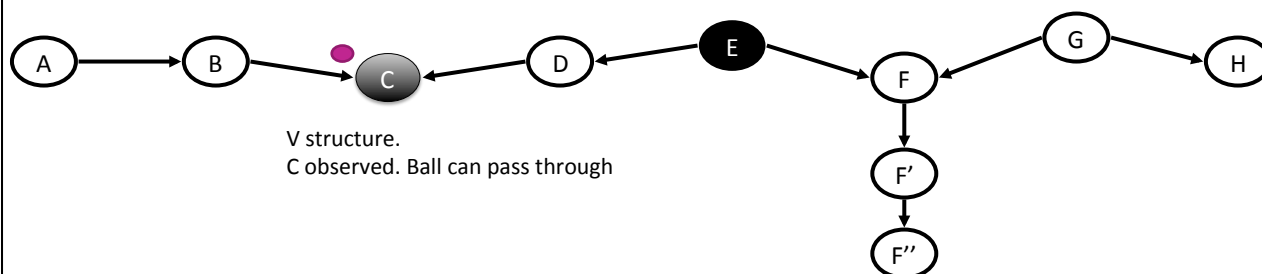
57

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



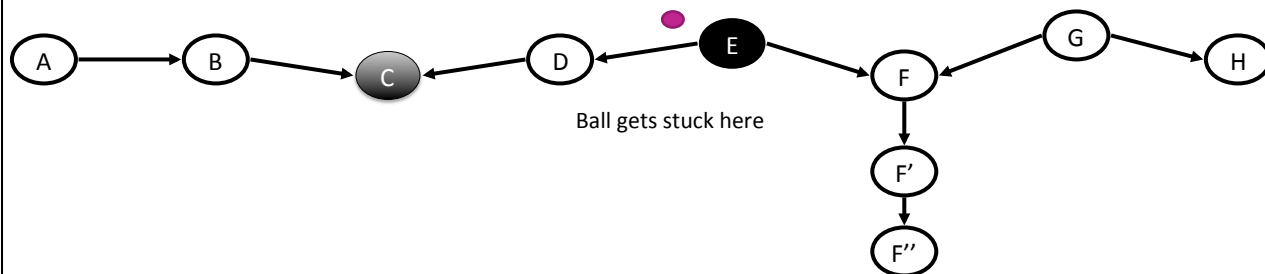
58

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



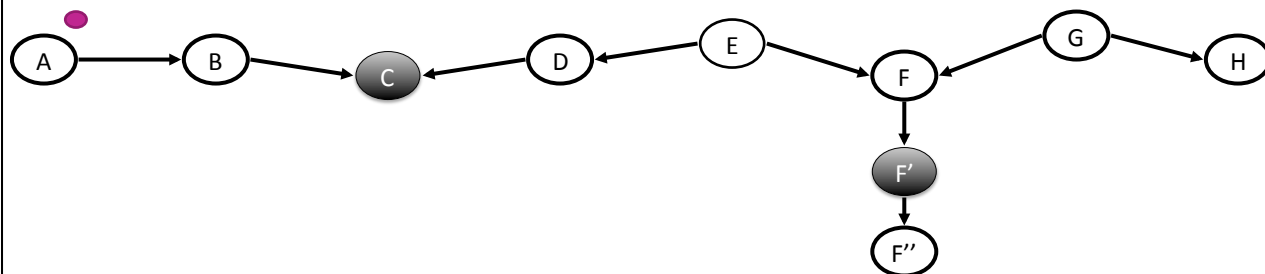
59

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



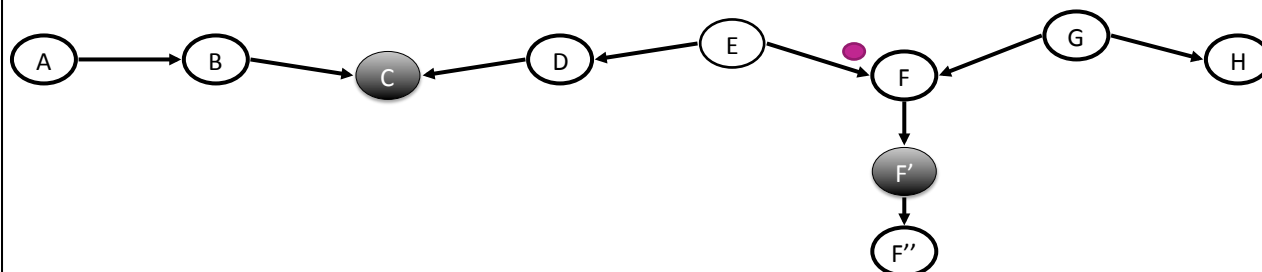
60

©2017 Emily Fox

CSE 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

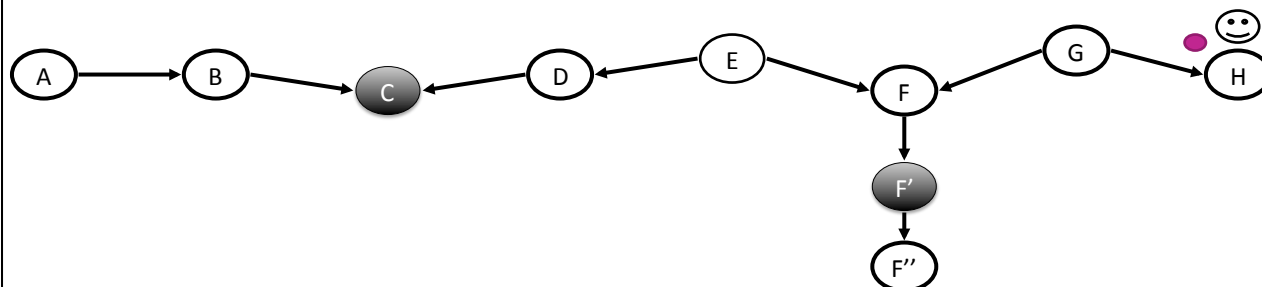
61

©2017 Emily Fox

CSE 446: Machine Learning

## Bayes ball example

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



62

©2017 Emily Fox

CSE 446: Machine Learning