

CS440/ECE448: Intro to Artificial Intelligence

# **Lecture 20**

## **More on learning graphical models**

Prof. Julia Hockenmaier  
[juliahmr@illinois.edu](mailto:juliahmr@illinois.edu)

<http://cs.illinois.edu/fa11/cs440>

# Bayes Nets

A Bayes Net defines a **joint distribution**  $P(X_1 \dots X_n)$  over a set of random variables  $X_1 \dots X_n$

Using the **chain rule**, we can factor  $P(X_1 \dots X_n)$  into a **product of  $n$  conditional distributions**:

$$P(X_1 \dots X_n) = \prod_j P(X_i \mid X_1 \dots X_{i-1}).$$

A Bayes Net makes a number of (conditional) independence assumptions:

$$P(X_1 \dots X_n) =_{\text{def}} \prod_j P(X_i \mid \text{Parents}(X_i) \subseteq \{X_1 \dots X_{i-1}\})$$

# Learning Bayes Nets

**Parameter estimation:** Given some data  $D$  over a set of random variables  $\mathbf{X}$  and a Bayes Net (with empty CPTs) estimate the parameters (= fill in the CPTs) of the Bayes Net.

**Structure learning:** Given some data  $D$  over a set of random variables  $\mathbf{X}$ , find a Bayes Net (define its CPTs) and estimate its parameters.

(This is much harder... we won't deal with it here)

# Bayes Rule

$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

$P(h)$ : prior probability of hypothesis

$P(h | D)$ : posterior probability of hypothesis.

$P(D | h)$ : likelihood of data, given hypothesis

Prior  $\propto$  posterior  $\times$  likelihood

$$P(h | D) \propto P(D | h)P(h)$$

# Three kinds of estimation techniques

Bayes optimal: Marginalize out the hypotheses

$$P(X | \mathbf{D}) = \sum_i P(X | h_i)P(h_i | \mathbf{D})$$

MAP (maximum a posteriori):

Pick the hypothesis with the highest posterior

$$h_{MAP} = \operatorname{argmax}_h P(h | \mathbf{D})$$

ML (maximum likelihood):

Pick the hypothesis that assigns highest likelihood

$$h_{ML} = \operatorname{argmax}_h P(\mathbf{D} | h)$$

# Maximum likelihood learning

Given data  $\mathbf{D}$ , we want to find the parameters that maximize  $P(\mathbf{D} \mid \theta)$ .

We have a data set with  $N$  candies.  
 $c$  are cherry.  $l = (N-c)$ , are lime.

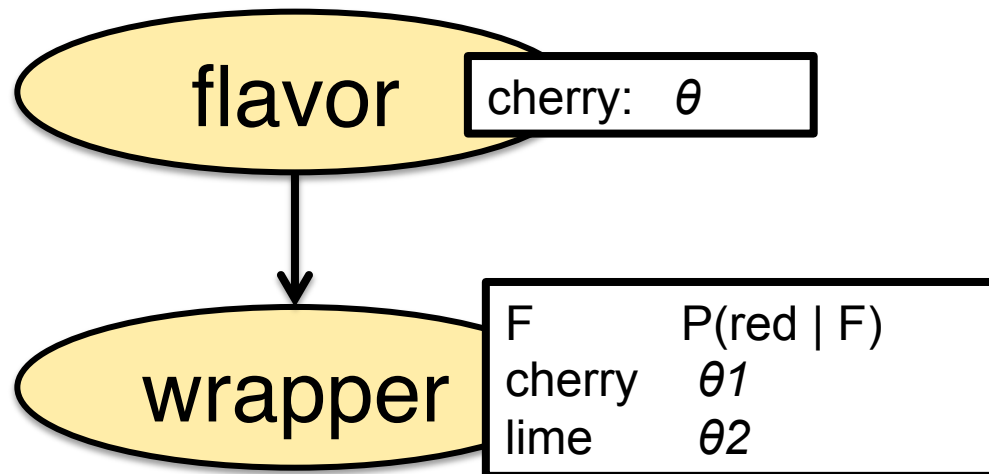
Parameter  $\theta$  = probability of cherry

Maximum likelihood estimate:  $\theta = c/N$

# A more complex model

Now the candy has two kinds of wrappers (red or green).

The wrapper is chosen probabilistically, depending on the flavor of the candy.



Out of  $N$  candies,  $c$  are cherry.  $r_c$  are cherry with a red wrapper,  $r_l$  are lime with a red wrapper

The **likelihood** of this data set:

$$P(d \mid \theta, \theta_1, \theta_2) = \theta^c (1-\theta)^{N-c} \theta_1^{r_c} (1-\theta_1)^{c-r_c} \theta_2^{r_l} (1-\theta_2)^{(N-c)-r_l}$$

The **log likelihood** of this data set:

$$\begin{aligned} L(d \mid \theta, \theta_1, \theta_2) = & [c \log \theta + (N-c) \log(1-\theta)] \\ & + [r_c \log \theta_1 + (c-r_c) \log(1-\theta_1)] \\ & + [r_l \log \theta_2 + (N-c-r_l) \log(1-\theta_2)] \end{aligned}$$

The **ML parameter estimates**:

$$\theta = c/N \quad \theta_1 = r_c/c \quad \theta_2 = r_l/(N-c)$$



# Medical diagnosis

Patients see a doctor and complain about a number of symptoms (headache, 100F fever, ...).

What is the most likely disease  $d_i$ , given the set of symptoms  $S$  the patient has?

$$\arg \max_{d_i} P(d_i | \bar{S})$$

# The Naïve Bayes classifier

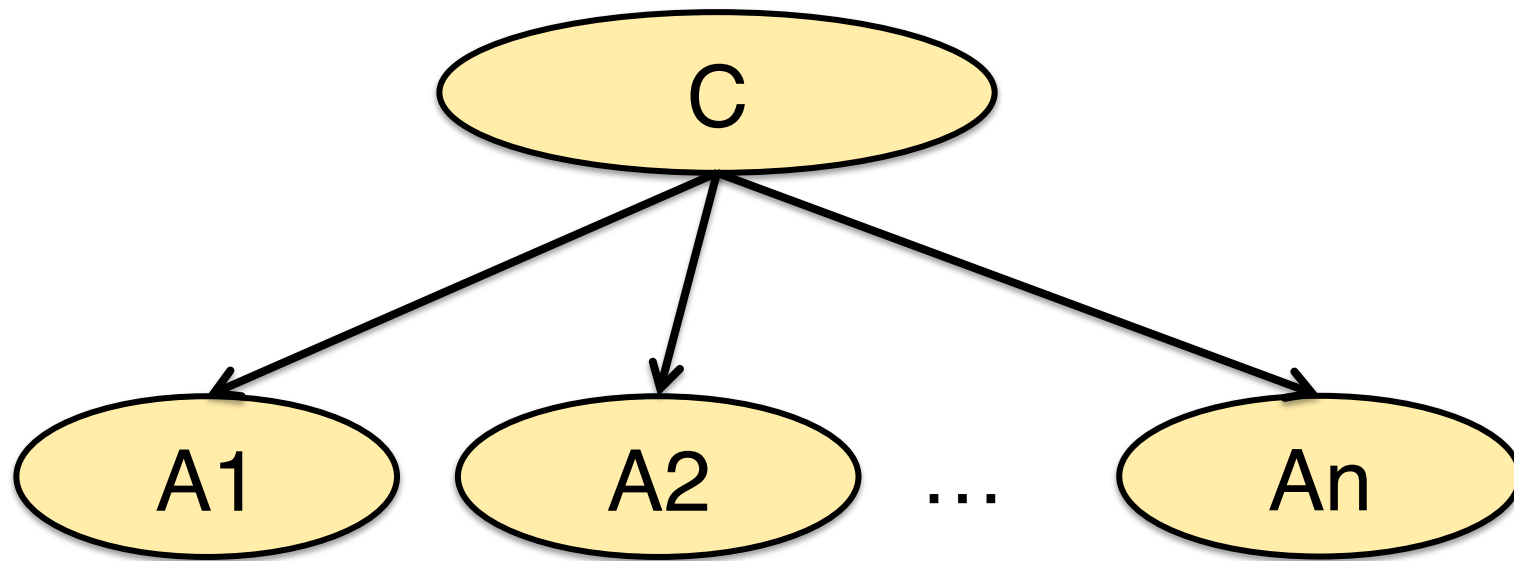
Assume the items in your data set have a number of attribute  $A_1 \dots A_n$ .

Each item also belongs to one of a number of given classes  $C_1 \dots C_k$ .

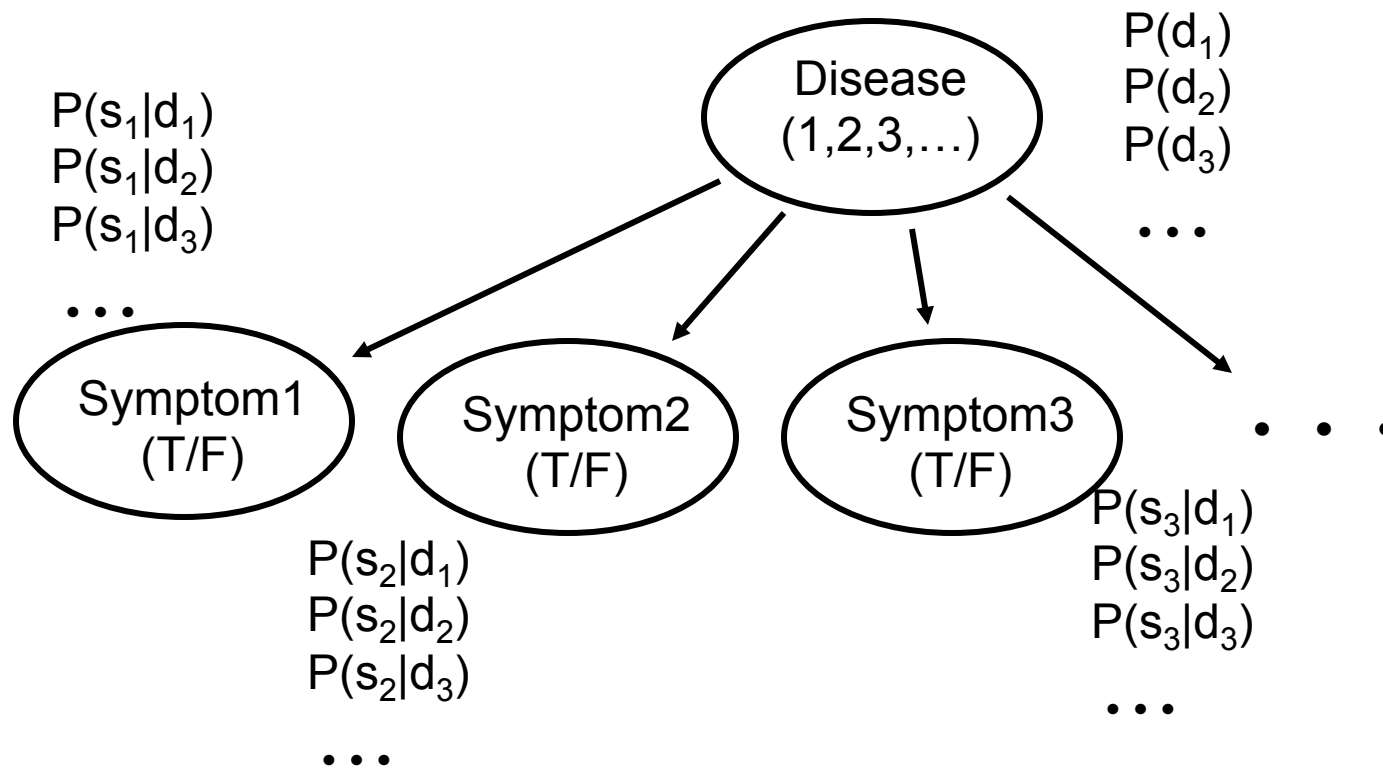
Which attributes an item has depends on its class.

If you only observe the attributes of an item, can you predict the class?

# The Naïve Bayes classifier



# Naïve Bayes



# Naïve Bayes

$$\begin{aligned}\operatorname{argmax}_C P(C | A_1 \dots A_n) &= \\ &= \operatorname{argmax}_C P(A_1 \dots A_n | C) P(C) \\ &= \operatorname{argmax}_C \prod_j P(A_j | C) P(C)\end{aligned}$$

We need to estimate:

- the multinomial  $P(C)$
- for each attribute  $A_j$  and class  $c$   $P(A_j | c)$

# Maximum likelihood estimation

If we have a set of training data where the class of each item is given:

- the multinomial  $P(C=c) = \text{freq}(c)/N$
- for each attribute  $A_j$  and class  $c$ :  
 $P(A_j = a | c) = \text{freq}(a, c) / \text{freq}(c)$

where

$\text{freq}(c)$  = the number of items in the training data that have class  $c$

$\text{freq}(a, c)$  = the number of items in the training data that have attribute  $a$  and class  $c$ .