

Lecture 15a: Statistical Learning

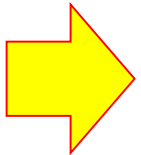
CSCI 360

Introduction to Artificial Intelligence

USC

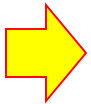
Here is where we are...

	3/1		Project 2 Out	
9	3/4 3/6	3/5 3/7	Quantifying Uncertainty Bayesian Networks	[Ch 13.1-13.6] [Ch 14.1-14.2]
10	3/11 3/13	3/12 3/14	(spring break, no class) (spring break, no class)	
11	3/18 3/20	3/19 3/21	Inference in Bayesian Networks Decision Theory	[Ch 14.3-14.4] [Ch 16.1-16.3 and 16.5]
	3/23		Project 2 Due	
12	3/25 3/27	3/26 3/28	<i>Advanced topics</i> (Chao traveling to National Science Foundation) <i>Advanced topics</i> (Chao traveling to National Science Foundation)	
	3/29		Homework 2 Out	
13	4/1 4/3	4/2 4/4	Markov Decision Processes Decision Tree Learning	[Ch 17.1-17.2] [Ch 18.1-18.3]
	4/5 4/5		Homework 2 Due Project 3 Out	
14	4/8 4/10	4/9 4/11	Perceptron Learning Neural Network Learning	[Ch 18.6] [Ch 18.7]
15	4/15 4/17	4/16 4/18	Statistical Learning Reinforcement Learning	[Ch 20.2.1-20.2.2] [Ch 21.1-21.2]
16	4/22 4/24	4/23 4/25	Artificial Intelligence Ethics Wrap-Up and Final Review	
	4/26		Project 3 Due	
	5/3	5/2	Final Exam (2pm-4pm)	



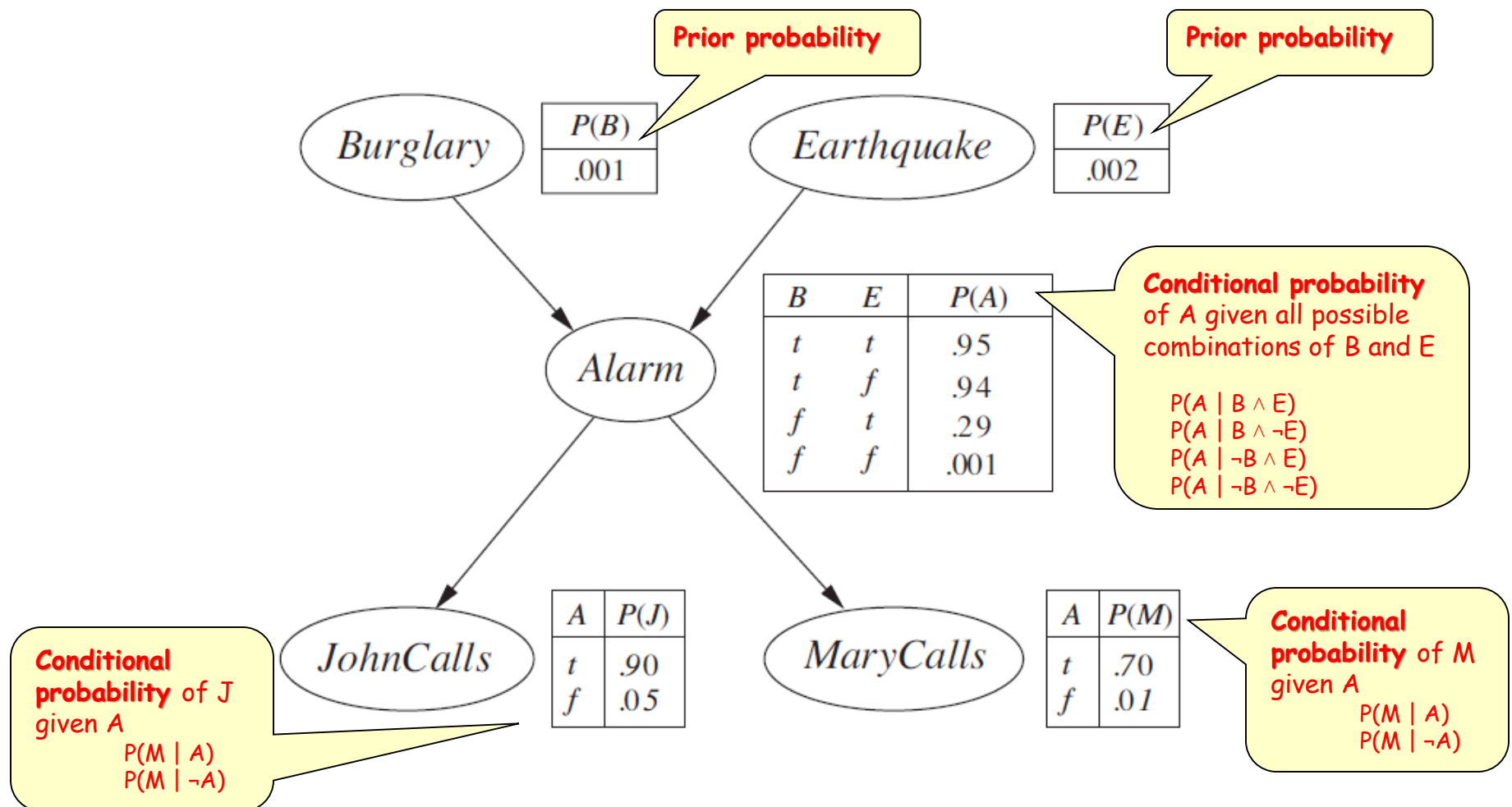
Outline

- What is AI?
- Part I: Search
- Part II: Logical reasoning
- Part III: Probabilistic reasoning
- **Part IV: Machine learning**
 - Decision Tree Learning
 - Perceptron Learning
 - Neural Network Learning
 - **Statistical Learning**
 - Reinforcement Learning



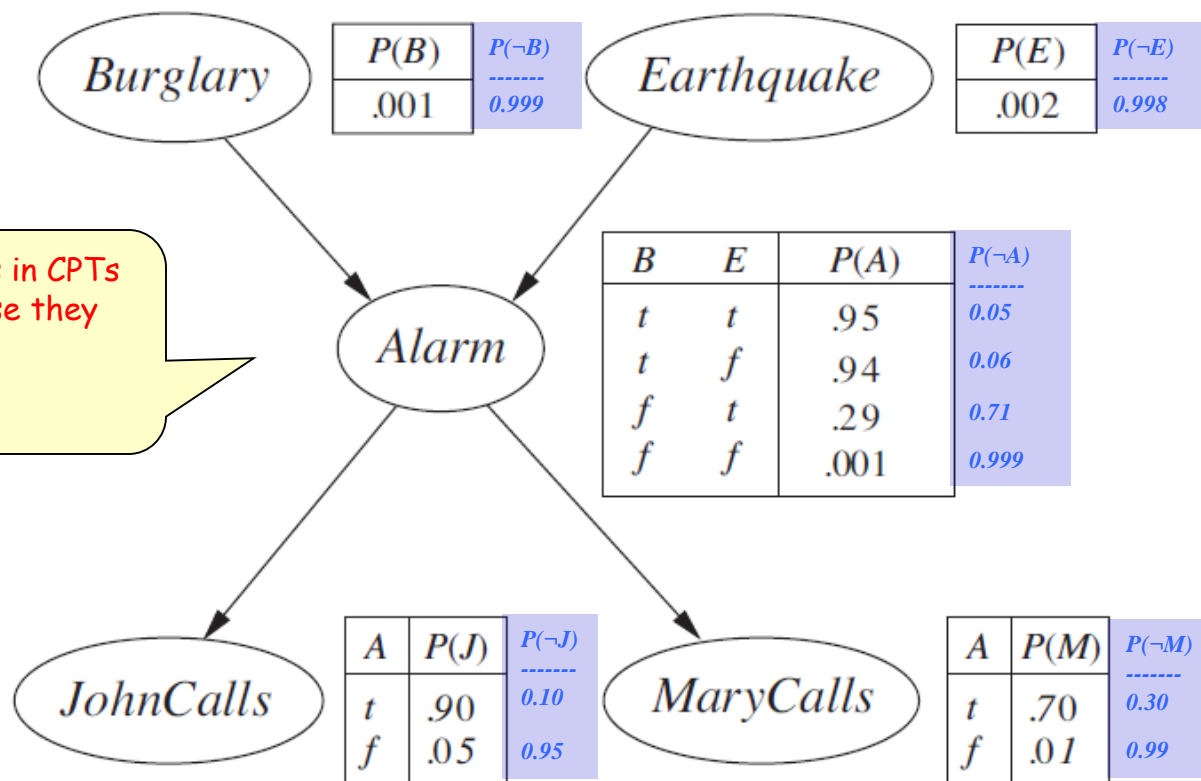
Recap: Bayesian networks (example)

- Both the **topology** and the **conditional probability tables** (CPTs)



Recap: *Bayesian networks (semantics)*

- Both the **topology** and the **conditional probability tables (CPTs)**

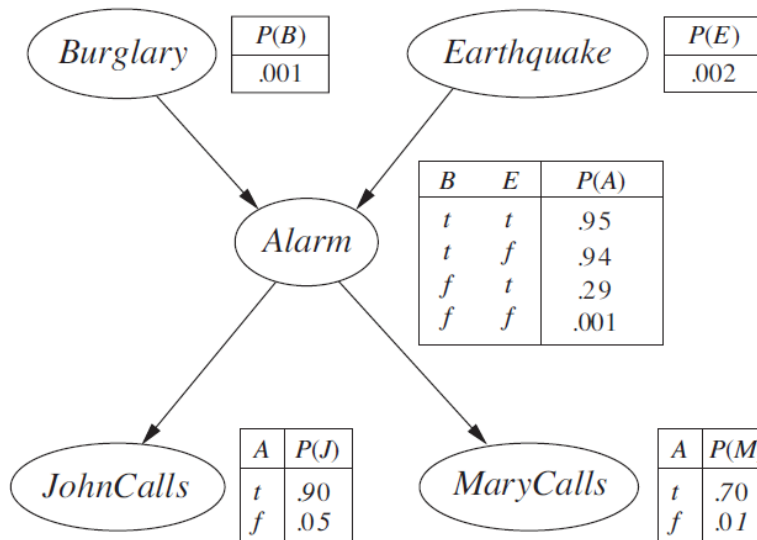


Half of the entries in CPTs are omitted because they can be inferred

Recap: Computing the joint distribution

- Each entry $P(x_1, \dots, x_n)$ in the full joint distribution, which is the abbreviation of $P(X_1=x_1 \wedge \dots \wedge X_n=x_n)$ is the product of the elements of the CPTs defined as follows:

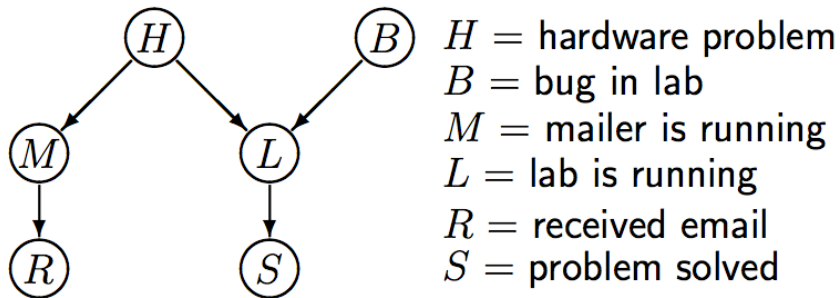
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$



$$\begin{aligned} P(j, m, a, b, e) \\ &= P(j/a) P(m/a) P(a/b, e) \mathbf{P(b)} \mathbf{P(e)} \\ &= 0.90 * 0.70 * 0.95 * \mathbf{0.001} * \mathbf{0.002} \\ &= 0.000001197 \end{aligned}$$

$$P(\neg m, j, \neg a, \neg e, b) = ?$$

Recap: *Bayesian inference*



Each node needs a probability table. Size of table depends on number of parents.

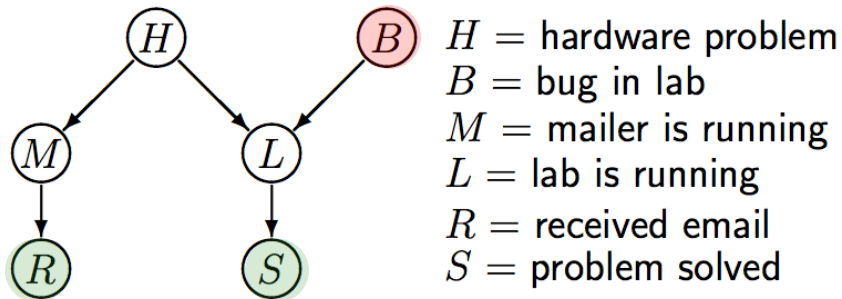
$\mathbf{P}(H)$	
<i>True</i>	<i>False</i>
0.01	0.99

H	$\mathbf{P}(M \mid H)$	
	<i>True</i>	<i>False</i>
<i>True</i>	0.1	0.9
<i>False</i>	0.99	0.01

H	B	$\mathbf{P}(L \mid H, B)$	
		<i>True</i>	<i>False</i>
<i>True</i>	<i>True</i>	0.01	0.99
<i>True</i>	<i>False</i>	0.1	0.9
<i>False</i>	<i>True</i>	0.02	0.98
<i>False</i>	<i>False</i>	1.0	0.0

..., etc.

Recap: *Bayesian inference*



- Compute $\mathbf{P}(B \mid \neg R, S)$

Each node needs a probability table. Size of table depends on number of parents.

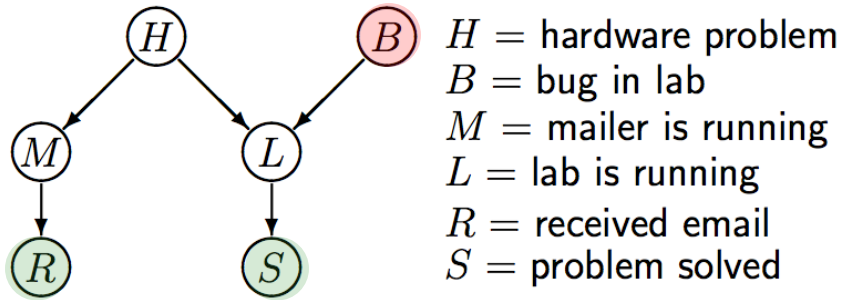
$\mathbf{P}(H)$	
<i>True</i>	<i>False</i>
0.01	0.99

	$\mathbf{P}(M \mid H)$	
H	<i>True</i>	<i>False</i>
<i>True</i>	0.1	0.9
<i>False</i>	0.99	0.01

		$\mathbf{P}(L \mid H, B)$	
H	B	<i>True</i>	<i>False</i>
<i>True</i>	<i>True</i>	0.01	0.99
<i>True</i>	<i>False</i>	0.1	0.9
<i>False</i>	<i>True</i>	0.02	0.98
<i>False</i>	<i>False</i>	1.0	0.0

..., etc.

Recap: *Bayesian inference*

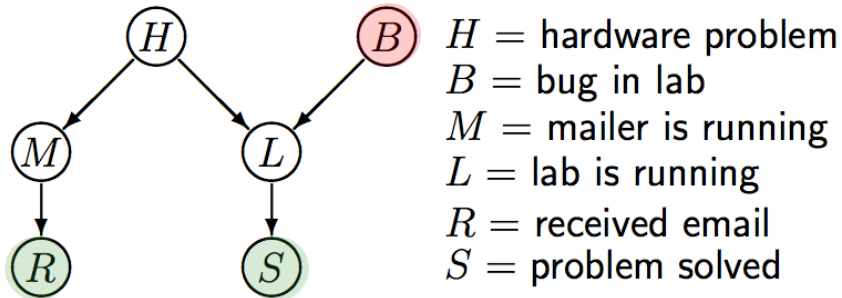


- Compute $\mathbf{P}(B \mid \neg R, S)$

1. Apply the conditional probability rule.

$$P(B \mid \neg R, S) = \frac{P(B, \neg R, S)}{P(\neg R, S)}$$

Recap: Bayesian inference



- Compute $\mathbf{P}(B \mid \neg R, S)$

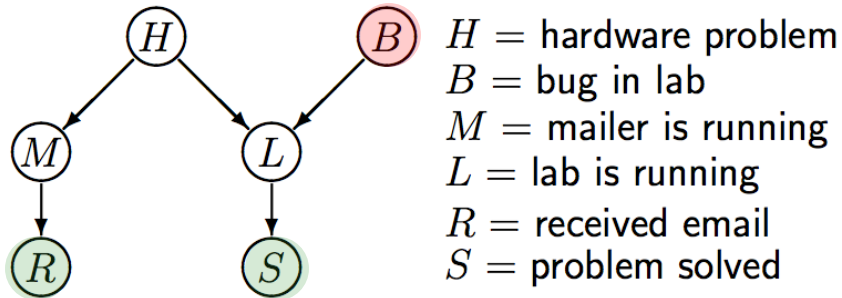
1. Apply the conditional probability rule.

$$P(B \mid \neg R, S) = \frac{P(B, \neg R, S)}{P(\neg R, S)}$$

2. Apply the marginal distribution rule to the unknown vertices. $P(B, \neg R, S)$ has 3 unknown vertices with $2^3 = 8$ possible value assignments.

$$\begin{aligned} P(B, \neg R, S) &= P(B, \neg R, S, H, M, L) \\ &\quad + P(B, \neg R, S, H, M, \neg L) \\ &\quad + P(B, \neg R, S, H, \neg M, L) \\ &\quad + P(B, \neg R, S, H, \neg M, \neg L) \\ &\quad + P(B, \neg R, S, \neg H, M, L) \\ &\quad + P(B, \neg R, S, \neg H, M, \neg L) \\ &\quad + P(B, \neg R, S, \neg H, \neg M, L) \\ &\quad + P(B, \neg R, S, \neg H, \neg M, \neg L) \end{aligned}$$

Recap: Bayesian inference



- Compute $\mathbf{P}(B \mid \neg R, S)$

1. Apply the conditional probability rule.

$$P(B \mid \neg R, S) = \frac{P(B, \neg R, S)}{P(\neg R, S)}$$

2. Apply the marginal distribution rule to the unknown vertices. $P(B, \neg R, S)$ has 3 unknown vertices with $2^3 = 8$ possible value assignments.

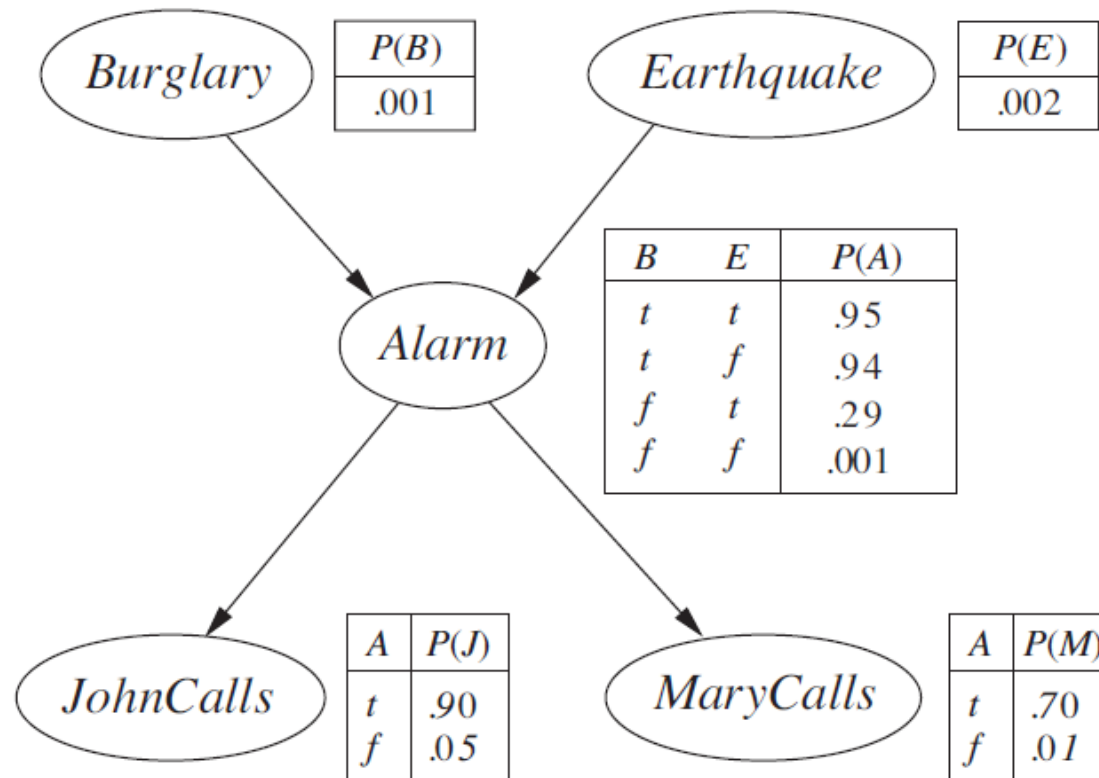
3. Apply joint distribution rule for Bayesian networks.

$$\begin{aligned}
 &P(B, \neg R, S, H, M, L) \\
 &= P(B) P(H) \\
 &\quad P(M \mid H) P(\neg R \mid M) \\
 &\quad P(L \mid H, B) P(S \mid L)
 \end{aligned}$$

$$\begin{aligned}
 &P(B, \neg R, S) \\
 &= P(B, \neg R, S, H, M, L) \\
 &\quad + P(B, \neg R, S, H, M, \neg L) \\
 &\quad + P(B, \neg R, S, H, \neg M, L) \\
 &\quad + P(B, \neg R, S, H, \neg M, \neg L) \\
 &\quad + P(B, \neg R, S, \neg H, M, L) \\
 &\quad + P(B, \neg R, S, \neg H, M, \neg L) \\
 &\quad + P(B, \neg R, S, \neg H, \neg M, L) \\
 &\quad + P(B, \neg R, S, \neg H, \neg M, \neg L)
 \end{aligned}$$

Question: *Where do the CPTs come from?*

- **Answer:** Learning probabilistic models from **data**



Outline of today's lecture

- **Statistical learning**
- Maximum-likelihood parameter learning
- Naïve Bayes models

Data and hypotheses

- Data (as the evidence)
 - Instantiations of random variables that describe the domain
- Hypotheses
 - Probabilistic theories of how the domain works
- Learning
 - Learn the **hypothesis** that best fit the **data**

Statistical learning example

- Candy in two flavors indistinguishable from the outside
 - cherry
 - lime
- Sold in large bags of five kinds:
 - h_1 : 100% cherry
 - h_2 : 75% cherry + 25% lime
 - h_3 : 50% cherry + 50% cherry
 - h_4 : 25% cherry + 75% lime
 - h_5 : 100% lime

Learning problem: Which bag ($H = h_1, \dots$, or $H = h_5$) is this ?

Statistical learning example

- Candy in two flavors indistinguishable from the outside
 - cherry
 - lime
- Sold in large bags of five kinds:
 - h_1 : 100% cherry
 - h_2 : 75% cherry + 25% lime
 - h_3 : 50% cherry + 50% cherry
 - h_4 : 25% cherry + 75% lime
 - h_5 : 100% lime

Learning problem: Which bag ($H = h_1, \dots$, or $H = h_5$), given favors of randomly drawn candies (D_1, \dots, D_N) ?

Bayesian learning

- Calculate the **probability of each hypothesis**, given the **data (\mathbf{d})**, and make the prediction on that basis

- $P(h_1 | \mathbf{d}) = \alpha P(\mathbf{d} | h_1) P(h_1)$

- $P(h_2 | \mathbf{d}) = \alpha P(\mathbf{d} | h_2) P(h_2)$

- $P(h_3 | \mathbf{d}) = \alpha P(\mathbf{d} | h_3) P(h_3)$

- $P(h_4 | \mathbf{d}) = \alpha P(\mathbf{d} | h_4) P(h_4)$

- $P(h_5 | \mathbf{d}) = \alpha P(\mathbf{d} | h_5) P(h_5)$

Likelihood of the data
under each hypothesis
 $P(\mathbf{d} | h_i)$

Hypothesis prior
 $P(h_i)$

Bayesian learning

- Calculate the **probability of each hypothesis**, given the **data (d)**, and make the prediction on that basis

$$\begin{aligned} - P(h_1 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_1) P(h_1) = \alpha P(\mathbf{d} | h_1) * 0.1 \\ - P(h_2 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_2) P(h_2) = \alpha P(\mathbf{d} | h_2) * 0.2 \\ - P(h_3 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_3) P(h_3) = \alpha P(\mathbf{d} | h_3) * 0.4 \\ - P(h_4 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_4) P(h_4) = \alpha P(\mathbf{d} | h_4) * 0.2 \\ - P(h_5 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_5) P(h_5) = \alpha P(\mathbf{d} | h_5) * 0.1 \end{aligned}$$

Likelihood of the data
under each hypothesis
 $P(\mathbf{d} | h_i)$

Hypothesis prior
 $P(h_i)$

Bayesian learning

- Calculate the **probability of each hypothesis**, given the **data** ($\mathbf{d} = \{d_1, \dots, d_5\}$), and make prediction on that basis

- $P(h_1 | \mathbf{d}) = \alpha P(\mathbf{d} | h_1) P(h_1) = \alpha P(\mathbf{d} | h_1) * 0.1$
- $P(h_2 | \mathbf{d}) = \alpha P(\mathbf{d} | h_2) P(h_2) = \alpha P(\mathbf{d} | h_2) * 0.2$
- $P(h_3 | \mathbf{d}) = \alpha P(\mathbf{d} | h_3) P(h_3) = \alpha P(\mathbf{d} | h_3) * 0.4$
- $P(h_4 | \mathbf{d}) = \alpha P(\mathbf{d} | h_4) P(h_4) = \alpha P(\mathbf{d} | h_4) * 0.2$
- $P(h_5 | \mathbf{d}) = \alpha P(\mathbf{d} | h_5) P(h_5) = \alpha P(\mathbf{d} | h_5) * 0.1$

Likelihood of the data
under each hypothesis
 $P(\mathbf{d} | h_i)$

Hypothesis prior
 $P(h_i)$

$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$

Bayesian learning

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction on that basis

- $P(h_1 | \mathbf{d}) = \alpha P(\mathbf{d} | h_1) P(h_1) = \alpha P(\mathbf{d} | h_1) * 0.1 = \alpha * 0 * 0.1$
- $P(h_2 | \mathbf{d}) = \alpha P(\mathbf{d} | h_2) P(h_2) = \alpha P(\mathbf{d} | h_2) * 0.2 = \alpha * 0.25^{10} * 0.2$
- $P(h_3 | \mathbf{d}) = \alpha P(\mathbf{d} | h_3) P(h_3) = \alpha P(\mathbf{d} | h_3) * 0.4 = \alpha * 0.50^{10} * 0.4$
- $P(h_4 | \mathbf{d}) = \alpha P(\mathbf{d} | h_4) P(h_4) = \alpha P(\mathbf{d} | h_4) * 0.2 = \alpha * 0.75^{10} * 0.2$
- $P(h_5 | \mathbf{d}) = \alpha P(\mathbf{d} | h_5) P(h_5) = \alpha P(\mathbf{d} | h_5) * 0.1 = \alpha * 1 * 0.1$

Likelihood of the data
under each hypothesis
 $P(\mathbf{d} | h_i)$

Hypothesis prior
 $P(h_i)$

$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$


First 10 candies are
all "lime"


Bayesian learning *(in general)*

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction of **unknown quantity X**

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

Each hypothesis (h_i) determines a probability distribution over X


$$P(h_i | \mathbf{d}) = \alpha P(\mathbf{d} | h_i) P(h_i)$$


$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$

Bayesian learning *(in general)*

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction of **unknown quantity X**

- Step 1. Computing the probability of data, given each hypothesis

$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$

- Step 2. Computing the probability of each hypothesis

$$P(h_i | \mathbf{d}) = \alpha P(\mathbf{d} | h_i) P(h_i)$$

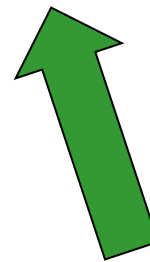
- Step 3. Computing the probability of X, given the data

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

Bayesian learning *(before any candy is revealed)*

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction on that basis

$$\begin{aligned} - P(h_1 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_1) * 0.1 = \alpha * 0.1 \\ - P(h_2 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_2) * 0.2 = \alpha * 0.2 \\ - P(h_3 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_3) * 0.4 = \alpha * \mathbf{0.4} \\ - P(h_4 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_4) * 0.2 = \alpha * 0.2 \\ - P(h_5 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_5) * 0.1 = \alpha * 0.1 \end{aligned}$$

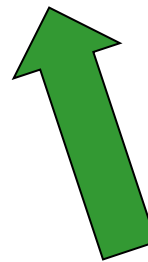


Know nothing about the candies in this bag

Bayesian learning *(after 1 candy is revealed)*

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction on that basis

- $P(h_1 | \mathbf{d}) = \alpha P(\mathbf{d} | h_1) * 0.1 = \alpha * 0 * 0.1 = \alpha * 0.00$
- $P(h_2 | \mathbf{d}) = \alpha P(\mathbf{d} | h_2) * 0.2 = \alpha * 0.25^1 * 0.2 = \alpha * 0.05$
- $P(h_3 | \mathbf{d}) = \alpha P(\mathbf{d} | h_3) * 0.4 = \alpha * 0.50^1 * 0.4 = \alpha * \mathbf{0.20}$
- $P(h_4 | \mathbf{d}) = \alpha P(\mathbf{d} | h_4) * 0.2 = \alpha * 0.75^1 * 0.2 = \alpha * 0.15$
- $P(h_5 | \mathbf{d}) = \alpha P(\mathbf{d} | h_5) * 0.1 = \alpha * 1 * 0.1 = \alpha * 0.10$

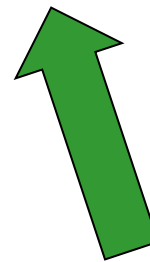


First candy is "lime"

Bayesian learning *(after 2 candies are revealed)*

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction on that basis

$$\begin{aligned} - P(h_1 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_1) * 0.1 = \alpha * 0 * 0.1 = \alpha * 0.0000 \\ - P(h_2 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_2) * 0.2 = \alpha * 0.25^2 * 0.2 = \alpha * 0.0125 \\ - P(h_3 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_3) * 0.4 = \alpha * 0.50^2 * 0.4 = \alpha * 0.1000 \\ - P(h_4 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_4) * 0.2 = \alpha * 0.75^2 * 0.2 = \alpha * \mathbf{0.1125} \\ - P(h_5 | \mathbf{d}) &= \alpha P(\mathbf{d} | h_5) * 0.1 = \alpha * 1 * 0.1 = \alpha * 0.1000 \end{aligned}$$

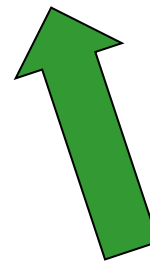


First 2 candies are
"lime"

Bayesian learning *(after 3 candies are revealed)*

- Calculate the **probability of each hypothesis**, given the **data**, and make the prediction on that basis

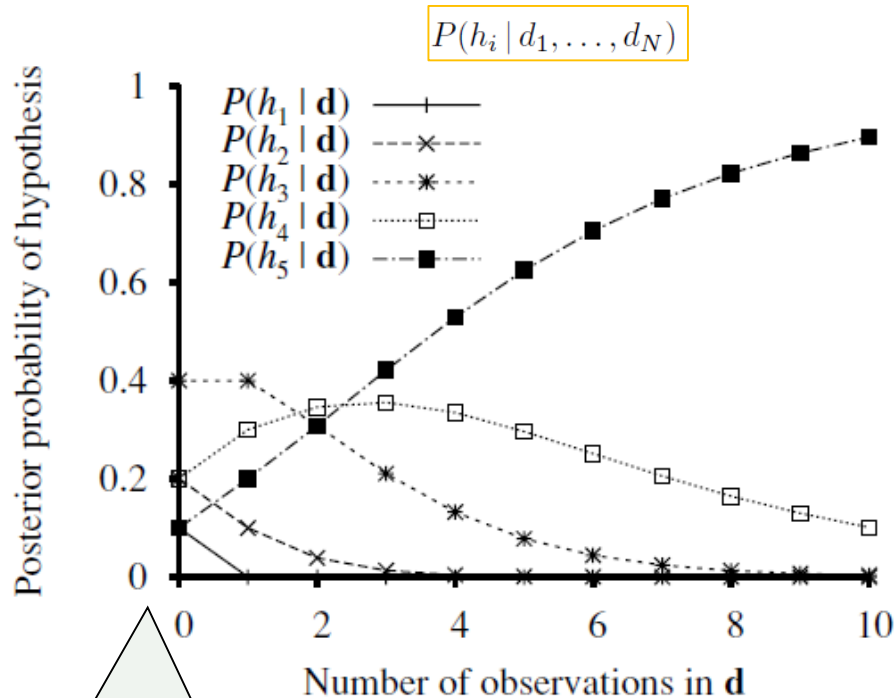
- $P(h_1 | \mathbf{d}) = \alpha P(\mathbf{d} | h_1) * 0.1 = \alpha * 0 * 0.1 = \alpha * 0.000000$
- $P(h_2 | \mathbf{d}) = \alpha P(\mathbf{d} | h_2) * 0.2 = \alpha * 0.25^2 * 0.2 = \alpha * 0.003125$
- $P(h_3 | \mathbf{d}) = \alpha P(\mathbf{d} | h_3) * 0.4 = \alpha * 0.50^3 * 0.4 = \alpha * 0.050000$
- $P(h_4 | \mathbf{d}) = \alpha P(\mathbf{d} | h_4) * 0.2 = \alpha * 0.75^3 * 0.2 = \alpha * 0.084375$
- $P(h_5 | \mathbf{d}) = \alpha P(\mathbf{d} | h_5) * 0.1 = \alpha * 1 * 0.1 = \alpha * 0.100000$



First 3 candies are
"lime"

Bayesian learning (results)

- Bayesian prediction **eventually** agrees with the **true hypothesis**

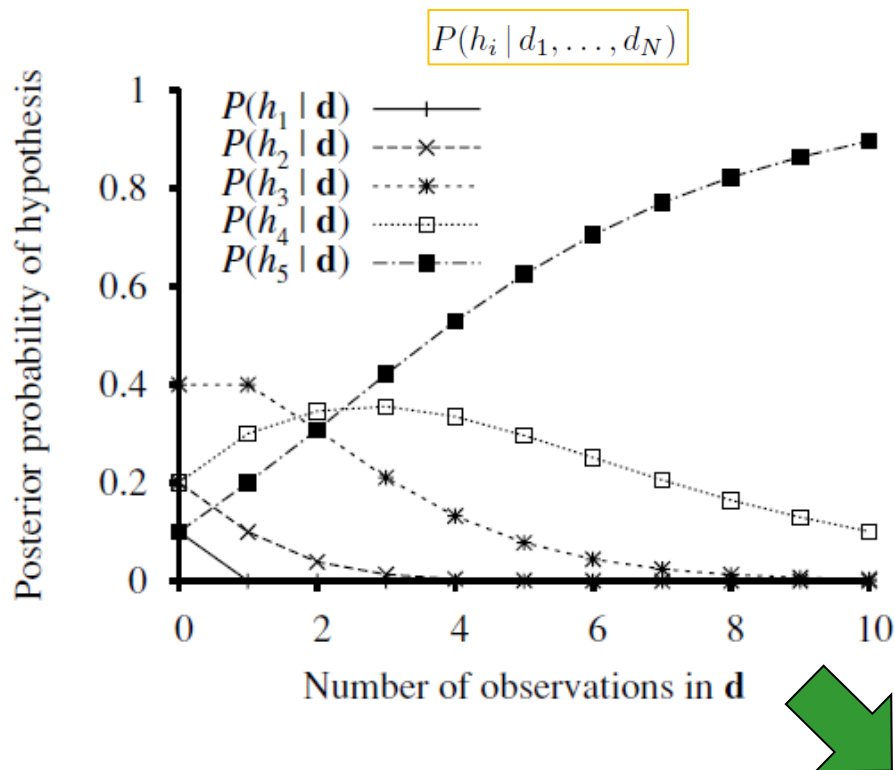


Eventually, the probability of the **true hypothesis** rises to the top

Initially, it is equal to the **prior** probability of the hypothesis

Bayesian learning (results)

- Bayesian prediction **eventually** agrees with the **true hypothesis**

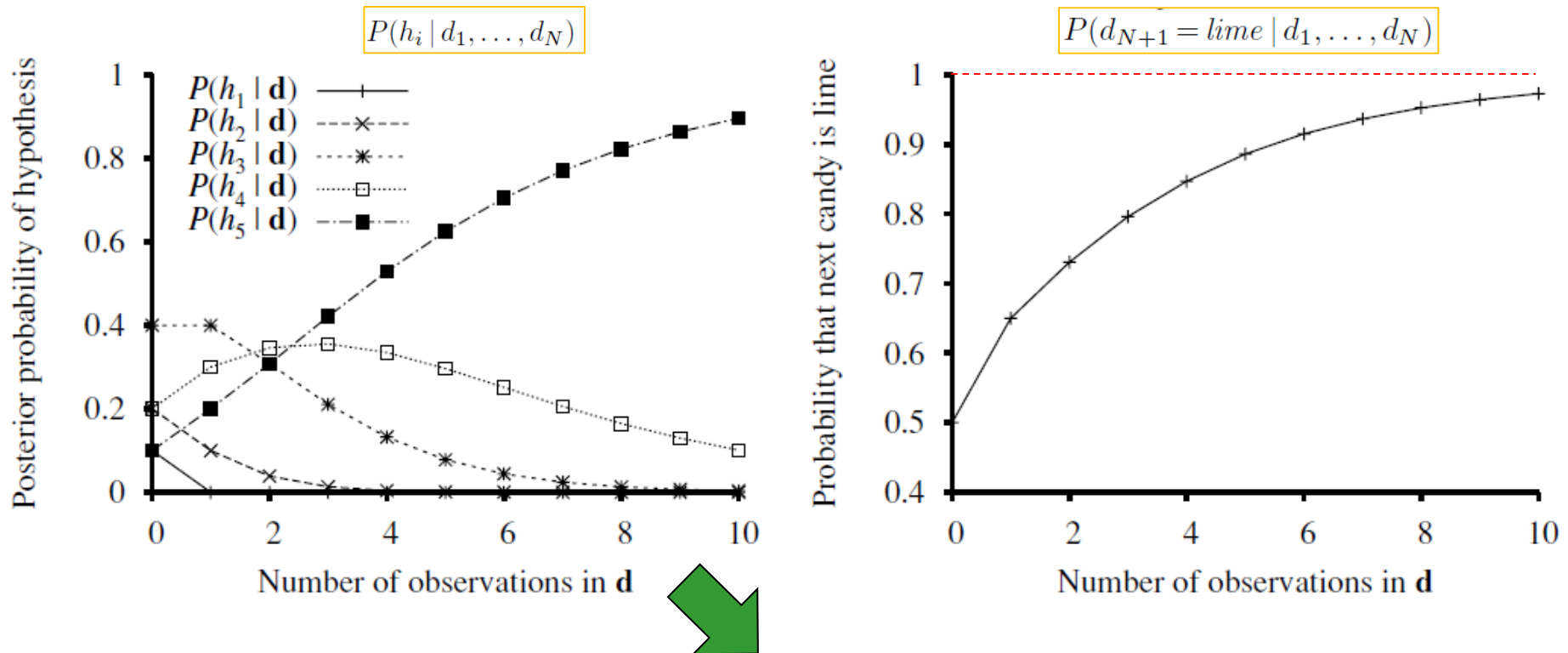


X is the quantity
to be predicted

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

Bayesian learning (results)

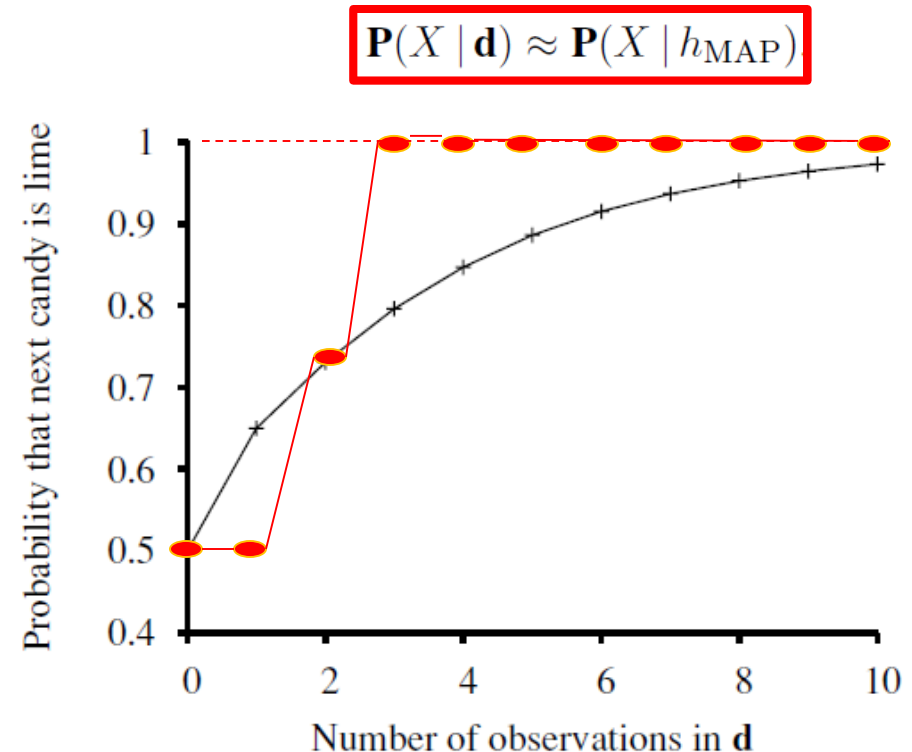
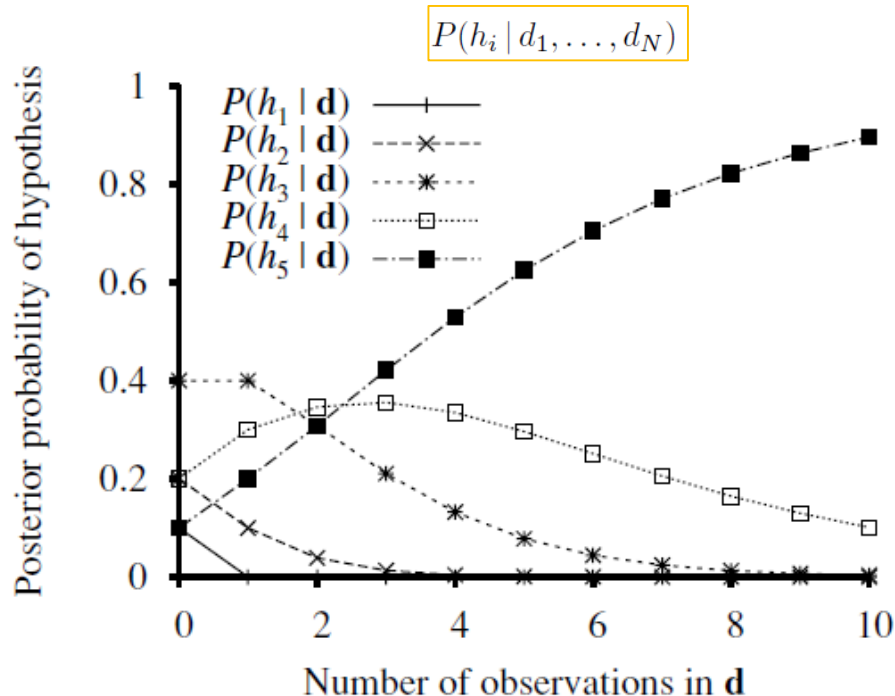
- Bayesian prediction **eventually** agrees with the **true hypothesis**



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

Maximum a posteriori (MAP)

- Make prediction based on a single, **most probable**, hypothesis



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

Bayesian learning (*MAP*)

- Approximating Bayesian prediction (**weighted sum**) by making decisions using the **most probable hypothesis**

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

$$P(X | \mathbf{d}) \approx P(X | h_{\text{MAP}})$$

h_{MAP} is the hypothesis
with largest $P(h_i | \mathbf{d})$

$$P(h_i | \mathbf{d}) = \alpha P(\mathbf{d} | h_i) P(h_i)$$

$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$

The role of prior – $P(h_i)$

- Bayesian (and MAP) learning uses **the prior** to penalize complexity
 - Complex hypothesis has a lower prior probability

$$\mathbf{P}(X | \mathbf{d}) \approx \mathbf{P}(X | h_{\text{MAP}})$$

h_{MAP} is the hypothesis with largest $P(h_i | \mathbf{d})$

$$P(h_i | \mathbf{d}) = \alpha P(\mathbf{d} | h_i) P(h_i)$$

$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$

Ockham's razor: find the simplest hypothesis that is consistent with the data

The role of prior – $P(h_i)$

- **Maximum-likelihood** hypothesis can be learned by using a **uniform prior**
 - Equivalent to maximizing $P(\mathbf{d} | h_i)$

$$\mathbf{P}(X | \mathbf{d}) \approx \mathbf{P}(X | h_{\text{MAP}})$$

h_{MAP} is the hypothesis with largest $P(h_i | \mathbf{d})$

$$P(h_i | \mathbf{d}) = \alpha P(\mathbf{d} | h_i) P(h_i)$$

$$P(\mathbf{d} | h_i) = \prod_j P(d_j | h_i)$$

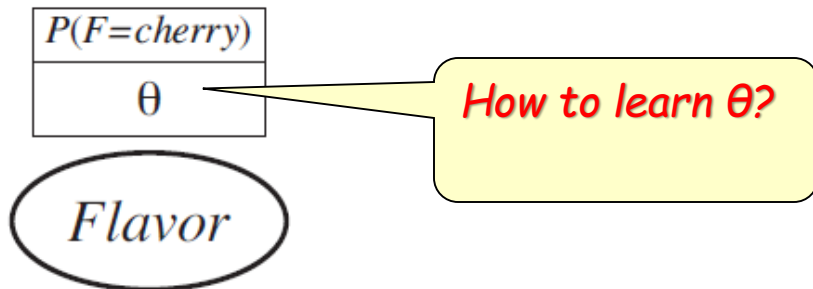
It provides a good approximation to Bayesian and MAP learning when the data set is large... but do not work well with a small data set

Outline of today's lecture

- Statistical learning
- **Maximum-likelihood parameter learning**
- Naïve Bayes models

Density estimation

- Learning a **probability model**, given **data** generated from that model
 - E.g., finding the **conditional probabilities** in a **Bayesian network** whose structure is fixed



Likelihood

- Let the proportions of “**cherry**” and “**lime**” candies be (θ) and ($1-\theta$), respectively.
- If we unwrap N candies, and find that there are
 - (c) **cherry** candies
 - ($l = N - c$) **lime** candies
- The likelihood of this data set (\mathbf{d}) is

$$P(\mathbf{d} | h_\theta) = \prod_{j=1}^N P(d_j | h_\theta) = \theta^c \cdot (1 - \theta)^l$$

$P(F=\text{cherry})$
θ

Flavor

Log likelihood

- Let the proportions of “**cherry**” and “**lime**” candies be (θ) and ($1-\theta$), respectively.
- If we unwrap N candies, and find that there are
 - (c) **cherry** candies
 - ($\ell = N - c$) **lime** candies
- The likelihood of this data set is

$$P(\mathbf{d} | h_\theta) = \prod_{j=1}^N P(d_j | h_\theta) = \theta^c \cdot (1 - \theta)^\ell$$



$$L(\mathbf{d} | h_\theta) = \log P(\mathbf{d} | h_\theta) = \sum_{j=1}^N \log P(d_j | h_\theta) = c \log \theta + \ell \log(1 - \theta)$$

Maximizing the log likelihood

- To find the maximum, compute the derivative of $L(d | h_\theta)$ and set it to zero.

$$L(\mathbf{d} | h_\theta) = \log P(\mathbf{d} | h_\theta) = \sum_{j=1}^N \log P(d_j | h_\theta) = c \log \theta + \ell \log(1 - \theta)$$

$$\frac{dL(\mathbf{d} | h_\theta)}{d\theta}$$

$$= \frac{c}{\theta} - \frac{\ell}{1 - \theta}$$

$$= 0$$

$$\Rightarrow \theta = \frac{c}{c + \ell} = \frac{c}{N}$$

Maximizing the log likelihood

- To find the maximum, compute the derivative of $L(\mathbf{d} | h_\theta)$ and set it to zero.

$$L(\mathbf{d} | h_\theta) = \log P(\mathbf{d} | h_\theta) = \sum_{j=1}^N \log P(d_j | h_\theta) = c \log \theta + \ell \log(1 - \theta)$$

$$\frac{dL(\mathbf{d} | h_\theta)}{d\theta}$$

$$= \frac{c}{\theta} - \frac{\ell}{1 - \theta}$$

$$= 0$$

The maximum-likelihood hypothesis (hML) says that the proportion (θ) of cherry candies is equal to the observed proportion

$$\Rightarrow \theta = \frac{c}{c + \ell} = \frac{c}{N}$$

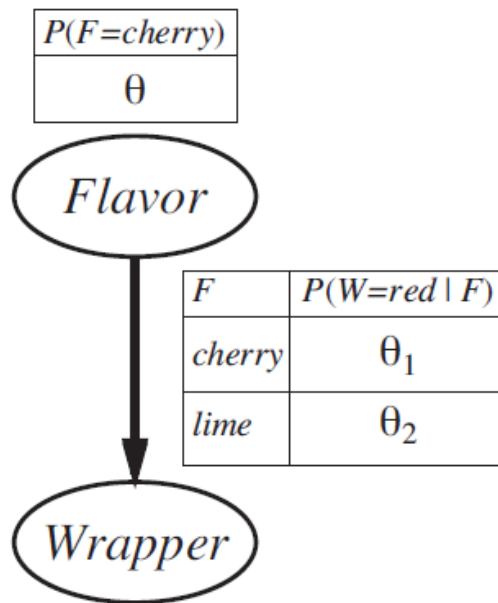
Obvious, but also comforting

General method *for maximum-likelihood parameter learning*

- Write down the **likelihood of the data** as a function of the parameters
- Compute the **derivative** of the **log likelihood** w.r.t. each parameter
- Find the parameter values such that the **derivatives are zero**

Another example

- Wrapper for each candy is chosen probabilistically based on the flavor
 - But the conditional distributions are unknown

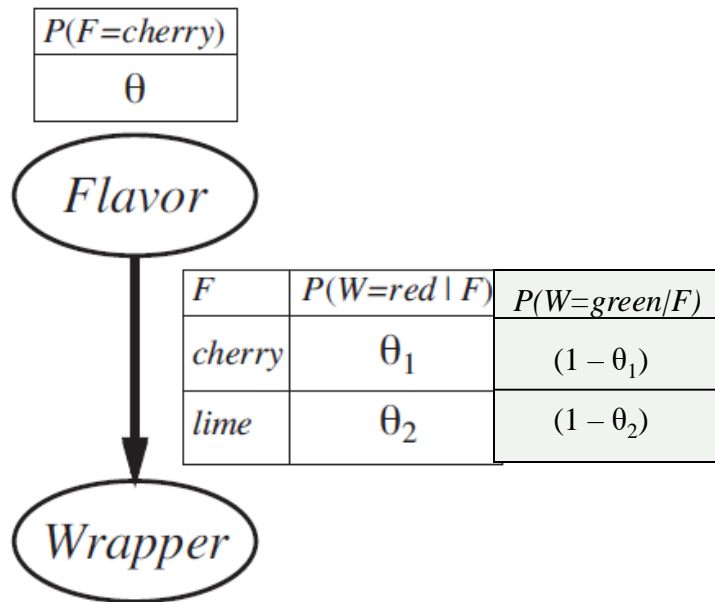


Question: What's the likelihood of seeing a **cherry** candy in a **green** wrapper?

Another example

- Wrapper for each candy is chosen probabilistically based the flavor
 - But the conditional distributions are unknown

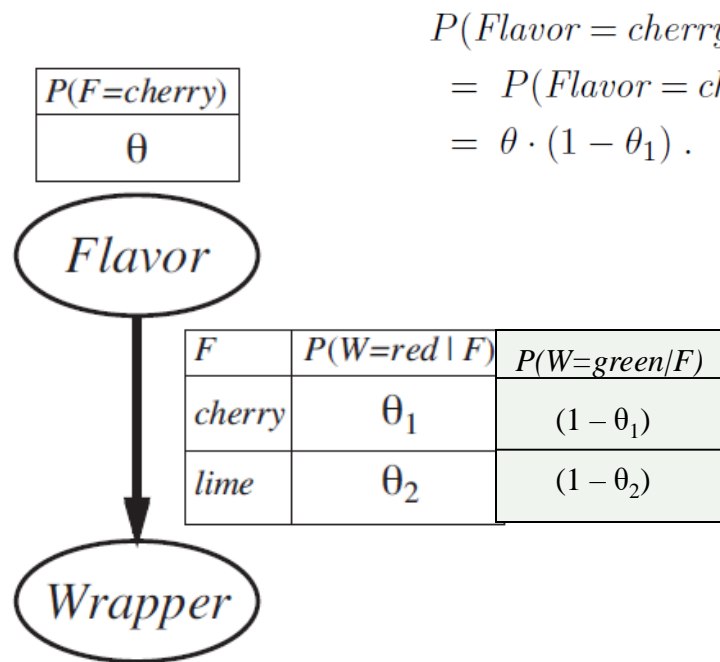
$$P(\text{Flavor} = \text{cherry}, \text{Wrapper} = \text{green} \mid h_{\theta, \theta_1, \theta_2})$$



Question: What's the likelihood of seeing a **cherry** candy in a **green** wrapper?

Another example

- Wrapper for each candy is chosen probabilistically based on the flavor
 - But the conditional distributions are unknown



$$\begin{aligned} P(\text{Flavor} = \text{cherry}, \text{Wrapper} = \text{green} \mid h_{\theta, \theta_1, \theta_2}) \\ &= P(\text{Flavor} = \text{cherry} \mid h_{\theta, \theta_1, \theta_2}) P(\text{Wrapper} = \text{green} \mid \text{Flavor} = \text{cherry}, h_{\theta, \theta_1, \theta_2}) \\ &= \theta \cdot (1 - \theta_1) . \end{aligned}$$

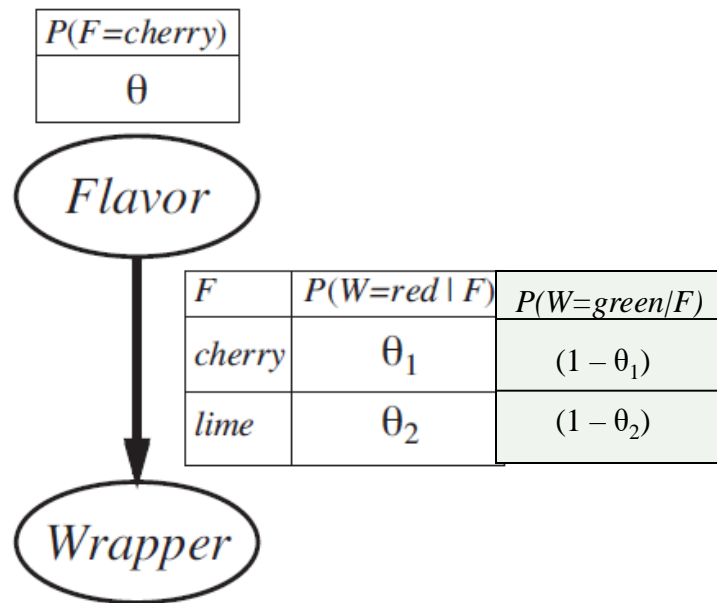
Question: What's the likelihood of seeing a **cherry** candy in a **green** wrapper?

Another example

- Data: Unwrapping **N** candies, of which **c** are cherries

$$P(\mathbf{d} \mid h_{\theta, \theta_1, \theta_2}) = \theta^c (1 - \theta)^\ell \cdot \theta_1^{r_c} (1 - \theta_1)^{g_c} \cdot \theta_2^{r_\ell} (1 - \theta_2)^{g_\ell}$$

$$L = [c \log \theta + \ell \log(1 - \theta)] + [r_c \log \theta_1 + g_c \log(1 - \theta_1)] + [r_\ell \log \theta_2 + g_\ell \log(1 - \theta_2)]$$



Another example

- Data: Unwrapping N candies, of which c are cherries

$$P(\mathbf{d} \mid h_{\theta, \theta_1, \theta_2}) = \theta^c (1 - \theta)^\ell \cdot \theta_1^{r_c} (1 - \theta_1)^{g_c} \cdot \theta_2^{r_\ell} (1 - \theta_2)^{g_\ell}$$

$$L = [c \log \theta + \ell \log(1 - \theta)] + [r_c \log \theta_1 + g_c \log(1 - \theta_1)] + [r_\ell \log \theta_2 + g_\ell \log(1 - \theta_2)]$$

- Derivatives

$$\begin{aligned}\frac{\partial L}{\partial \theta} &= \frac{c}{\theta} - \frac{\ell}{1 - \theta} = 0 \\ \frac{\partial L}{\partial \theta_1} &= \frac{r_c}{\theta_1} - \frac{g_c}{1 - \theta_1} = 0 \\ \frac{\partial L}{\partial \theta_2} &= \frac{r_\ell}{\theta_2} - \frac{g_\ell}{1 - \theta_2} = 0\end{aligned}$$

Another example

- Data: Unwrapping N candies, of which c are cherries

$$P(\mathbf{d} \mid h_{\theta, \theta_1, \theta_2}) = \theta^c (1 - \theta)^\ell \cdot \theta_1^{r_c} (1 - \theta_1)^{g_c} \cdot \theta_2^{r_\ell} (1 - \theta_2)^{g_\ell}$$

$$L = [c \log \theta + \ell \log(1 - \theta)] + [r_c \log \theta_1 + g_c \log(1 - \theta_1)] + [r_\ell \log \theta_2 + g_\ell \log(1 - \theta_2)]$$

- Derivatives

$$\begin{aligned} \frac{\partial L}{\partial \theta} &= \frac{c}{\theta} - \frac{\ell}{1-\theta} = 0 & \Rightarrow \theta &= \frac{c}{c+\ell} \\ \frac{\partial L}{\partial \theta_1} &= \frac{r_c}{\theta_1} - \frac{g_c}{1-\theta_1} = 0 & \Rightarrow \theta_1 &= \frac{r_c}{r_c+g_c} \\ \frac{\partial L}{\partial \theta_2} &= \frac{r_\ell}{\theta_2} - \frac{g_\ell}{1-\theta_2} = 0 & \Rightarrow \theta_2 &= \frac{r_\ell}{r_\ell+g_\ell} \end{aligned}$$

Another example

- Data: Unwrapping **N** candies, of which **c** are cherries

$$P(\mathbf{d} \mid h_{\theta, \theta_1, \theta_2}) = \theta^c (1 - \theta)^\ell \cdot \theta_1^{r_c} (1 - \theta_1)^{g_c} \cdot \theta_2^{r_\ell} (1 - \theta_2)^{g_\ell}$$

$$L = [c \log \theta + \ell \log(1 - \theta)] + [r_c \log \theta_1 + g_c \log(1 - \theta_1)] + [r_\ell \log \theta_2 + g_\ell \log(1 - \theta_2)]$$

- Derivatives

$$\begin{aligned}\frac{\partial L}{\partial \theta} &= \frac{c}{\theta} - \frac{\ell}{1-\theta} = 0 \\ \frac{\partial L}{\partial \theta_1} &= \frac{r_c}{\theta_1} - \frac{g_c}{1-\theta_1} = 0 \\ \frac{\partial L}{\partial \theta_2} &= \frac{r_\ell}{\theta_2} - \frac{g_\ell}{1-\theta_2} = 0\end{aligned}$$

$$\begin{aligned}\Rightarrow \theta &= \frac{c}{c+\ell} \\ \Rightarrow \theta_1 &= \frac{r_c}{r_c+g_c} \\ \Rightarrow \theta_2 &= \frac{r_\ell}{r_\ell+g_\ell}\end{aligned}$$

Same as before - obvious,
but also comforting

Proportion is equal to the
"observed" proportion

Another example

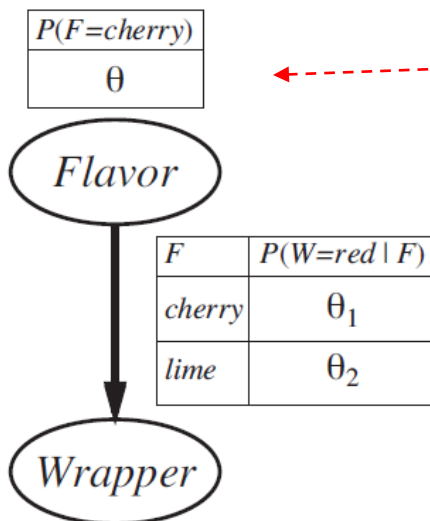
- Data: Unwrapping **N** candies, of which **c** are cherries

$$P(\mathbf{d} \mid h_{\theta, \theta_1, \theta_2}) = \theta^c (1 - \theta)^\ell \cdot \theta_1^{r_c} (1 - \theta_1)^{g_c} \cdot \theta_2^{r_\ell} (1 - \theta_2)^{g_\ell}$$

$$L = [c \log \theta + \ell \log(1 - \theta)] + [r_c \log \theta_1 + g_c \log(1 - \theta_1)] + [r_\ell \log \theta_2 + g_\ell \log(1 - \theta_2)]$$

Learning is **compositional**,
one for each parameter

Same as before - obvious,
but also comforting



$$\theta = \frac{c}{c + \ell}$$
$$\theta_1 = \frac{r_c}{r_c + g_c}$$
$$\theta_2 = \frac{r_\ell}{r_\ell + g_\ell}$$

Proportion is equal to the
"observed" proportion

Outline of today's lecture

- Statistical learning
- Maximum-likelihood parameter learning
- **Naïve Bayes models**

Recap: *Naïve Bayes model*

- A single **cause** directly influence a number of **effects**, all of which are conditionally independent, given the cause

$$\mathbf{P}(Cause, Effect_1, \dots, Effect_n) = \mathbf{P}(Cause) \prod_i \mathbf{P}(Effect_i \mid Cause)$$

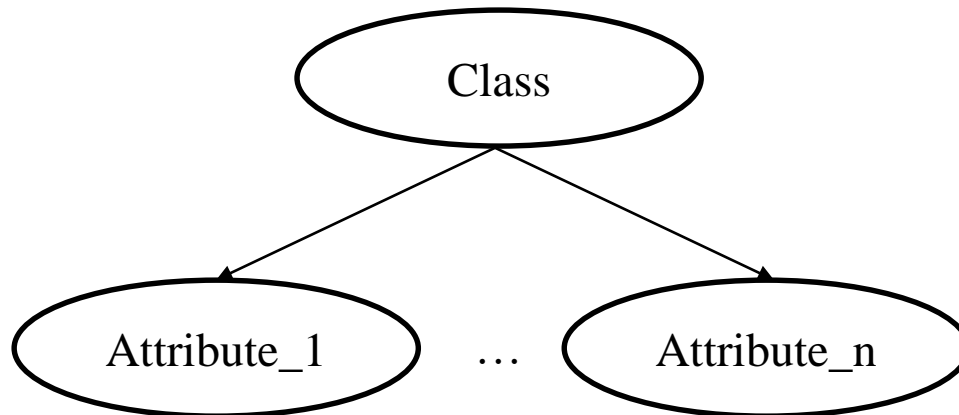
$$\begin{aligned} & \mathbf{P}(Toothache, Catch, Cavity) \\ &= \mathbf{P}(Toothache, Catch \mid Cavity) \mathbf{P}(Cavity) \\ &= \mathbf{P}(Toothache \mid Cavity) \mathbf{P}(Catch \mid Cavity) \mathbf{P}(Cavity) \end{aligned}$$

Naïve Bayesian learning

- **Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

- Using **maximum-likelihood estimates** to learn CDTs; that is, using “frequencies” to compute the “probabilities”



Naïve Bayesian learning

- Properties
 - Tolerant of noise in **attribute** and **class** values of examples
 - Learn quickly, even for large problems
- Early application
 - Email spam detector, where
 - *Attribute_i* = “how often does the *i*-th word in a dictionary appear in the email?”
 - *Class* = “is this email spam?”

How to use the learned model?

- **Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

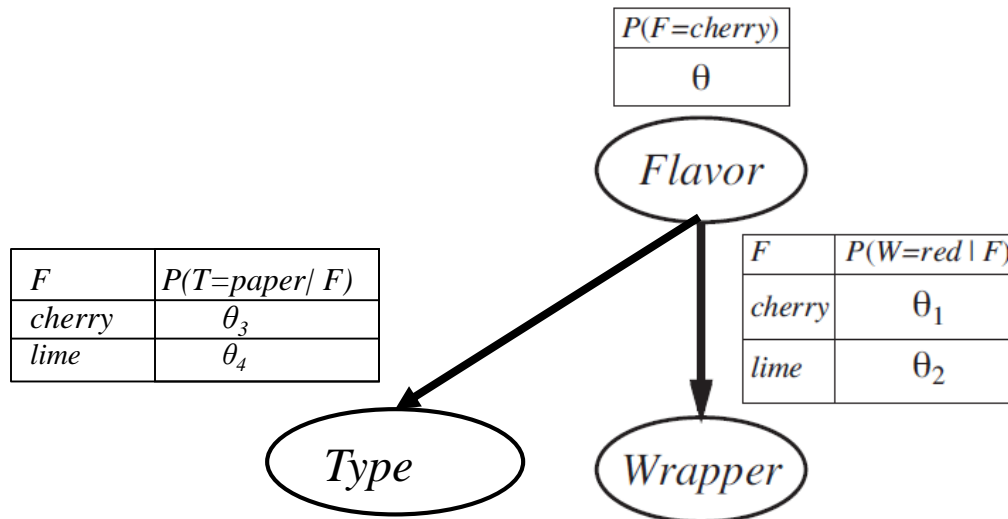
With observed attribute values x_1, \dots, x_n , what's the probability of each class C ?

How to use the learned model?

- **Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

With observed attribute values x_1, \dots, x_n , what's the probability of each class C ?



The **wrapper** color may be "red" or "green"

The material **type** may be "paper" or "plastic"

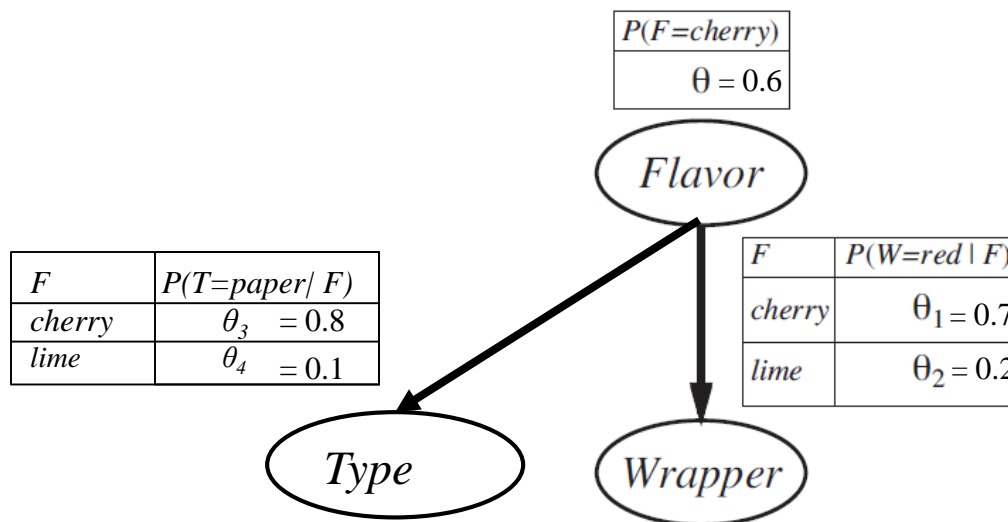
Example

- Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

$P(F = \text{cherry} \mid W = \text{red}, T = \text{paper}) = ?$

$P(F = \text{lime} \mid W = \text{red}, T = \text{paper}) = ?$



The **wrapper** color may be "red" or "green"

The material **type** may be "paper" or "plastic"

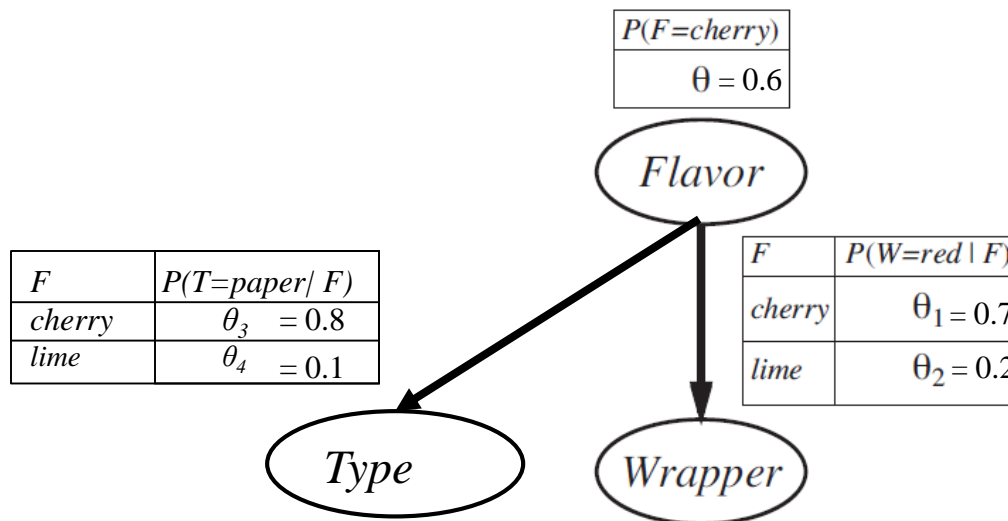
Example

- Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

$$\begin{aligned} P(F = \text{cherry} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{cherry}) * P(\text{red} \mid \text{cherry}) * P(\text{paper} \mid \text{cherry}) \\ &= \alpha * 0.6 * 0.7 * 0.8 = \alpha * 0.336 \end{aligned}$$

$$P(F = \text{lime} \mid W = \text{red}, T = \text{paper}) = ?$$



The **wrapper** color may be "red" or "green"

The material **type** may be "paper" or "plastic"

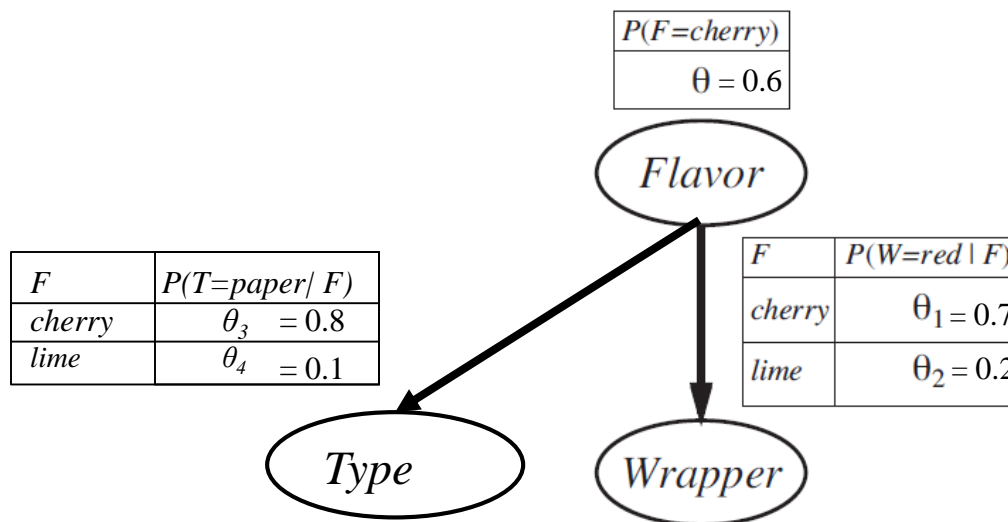
Example

- Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

$$\begin{aligned} P(F = \text{cherry} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{cherry}) * P(\text{red} \mid \text{cherry}) * P(\text{paper} \mid \text{cherry}) \\ &= \alpha * 0.6 * 0.7 * 0.8 = \alpha * 0.336 \end{aligned}$$

$$\begin{aligned} P(F = \text{lime} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{lime}) * P(\text{red} \mid \text{lime}) * P(\text{paper} \mid \text{lime}) \\ &= \alpha * 0.4 * 0.2 * 0.1 = \alpha * 0.008 \end{aligned}$$



The **wrapper** color may be "red" or "green"

The material **type** may be "paper" or "plastic"

Example

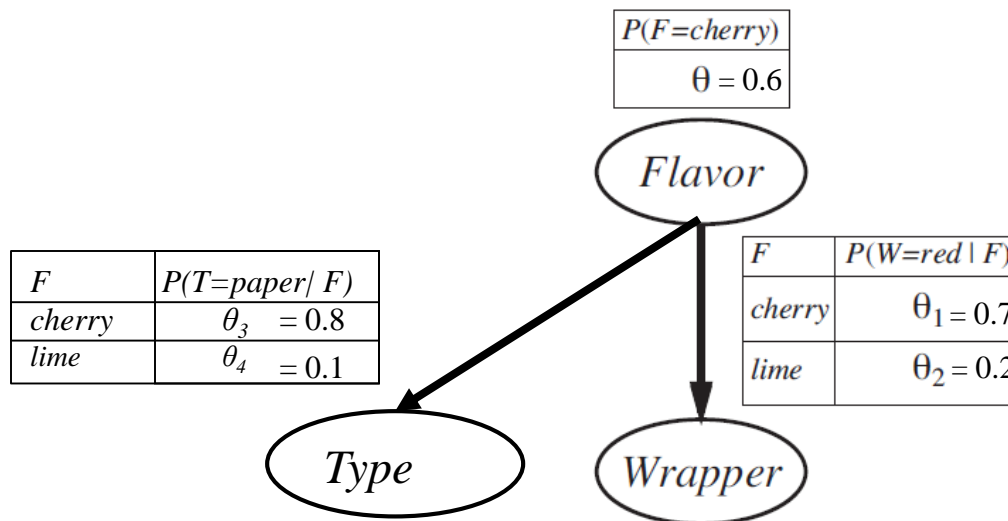
- Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

$$\begin{aligned} P(F = \text{cherry} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{cherry}) * P(\text{red} \mid \text{cherry}) * P(\text{paper} \mid \text{cherry}) \\ &= \alpha * 0.6 * 0.7 * 0.8 = \alpha * 0.336 \end{aligned}$$

$$\begin{aligned} P(F = \text{lime} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{lime}) * P(\text{red} \mid \text{lime}) * P(\text{paper} \mid \text{lime}) \\ &= \alpha * 0.4 * 0.2 * 0.1 = \alpha * 0.008 \end{aligned}$$

$$\alpha = 1 / (0.336 + 0.008) = 1 / 0.344$$



The **wrapper** color may be "red" or "green"

The material **type** may be "paper" or "plastic"

Example

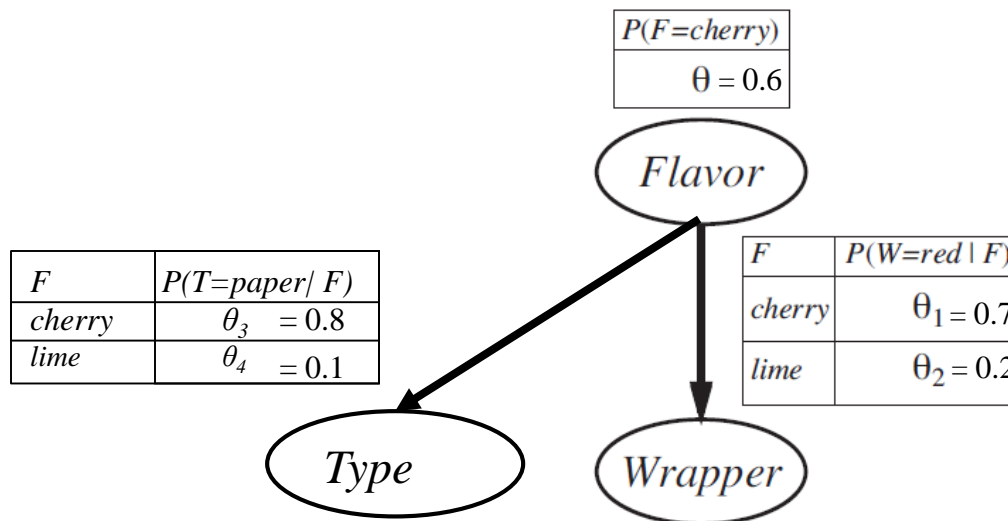
- Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

$$\begin{aligned} P(F = \text{cherry} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{cherry}) * P(\text{red} \mid \text{cherry}) * P(\text{paper} \mid \text{cherry}) \\ &= \alpha * 0.6 * 0.7 * 0.8 = \alpha * 0.336 = \mathbf{0.977 (97.7\%)} \end{aligned}$$

$$\begin{aligned} P(F = \text{lime} \mid W = \text{red}, T = \text{paper}) &= \alpha * P(\text{lime}) * P(\text{red} \mid \text{lime}) * P(\text{paper} \mid \text{lime}) \\ &= \alpha * 0.4 * 0.2 * 0.1 = \alpha * 0.008 = \mathbf{0.023 (2.3\%)} \end{aligned}$$

$$\alpha = 1 / (0.336 + 0.008) = 1 / 0.344$$



The **wrapper** color may be "red" or "green"

The material **type** may be "paper" or "plastic"

Quiz 13

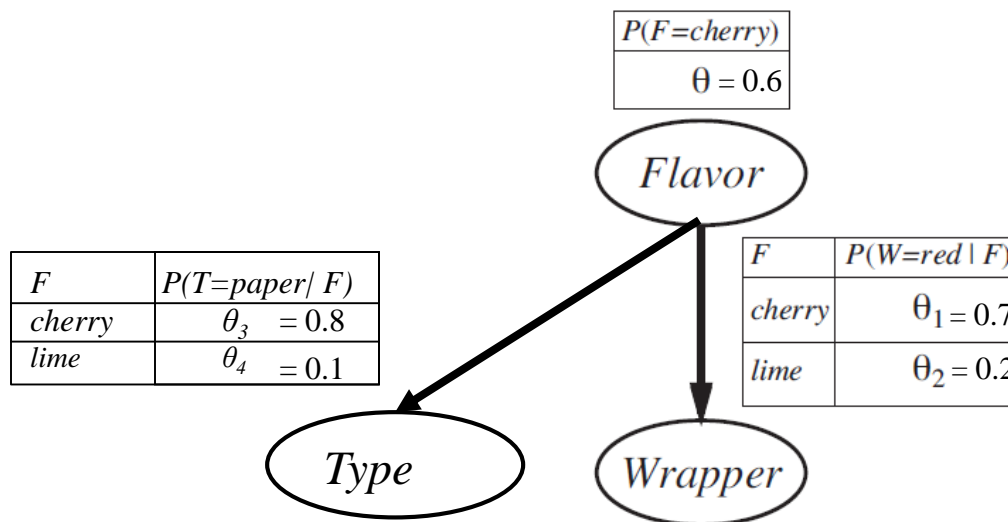
- Assumption:** Attributes (X_1, \dots, X_n) are conditionally independent of each other, given the class (C)

$$\mathbf{P}(C \mid x_1, \dots, x_n) = \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)$$

With the observed wrapper being “green” and “plastic”, what is the probability of the candy being “cherry” and “lime”, respectively?

$P(F = \text{cherry} \mid W = \text{green}, T = \text{plastic}) = ?$

$P(F = \text{lime} \mid W = \text{green}, T = \text{plastic}) = ?$



The **wrapper** color may be “red” or “green”

The material **type** may be “paper” or “plastic”