

# CS 1571 Introduction to AI

## Lecture ?

## Machine Learning of BBNs from data

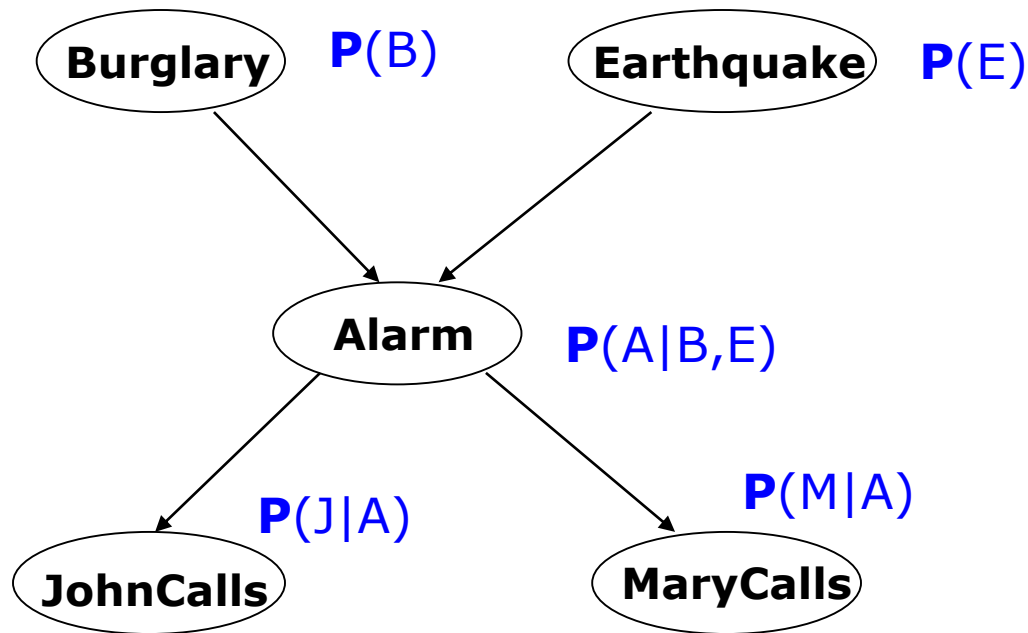
**Milos Hauskrecht**

[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)

5329 Sennott Square

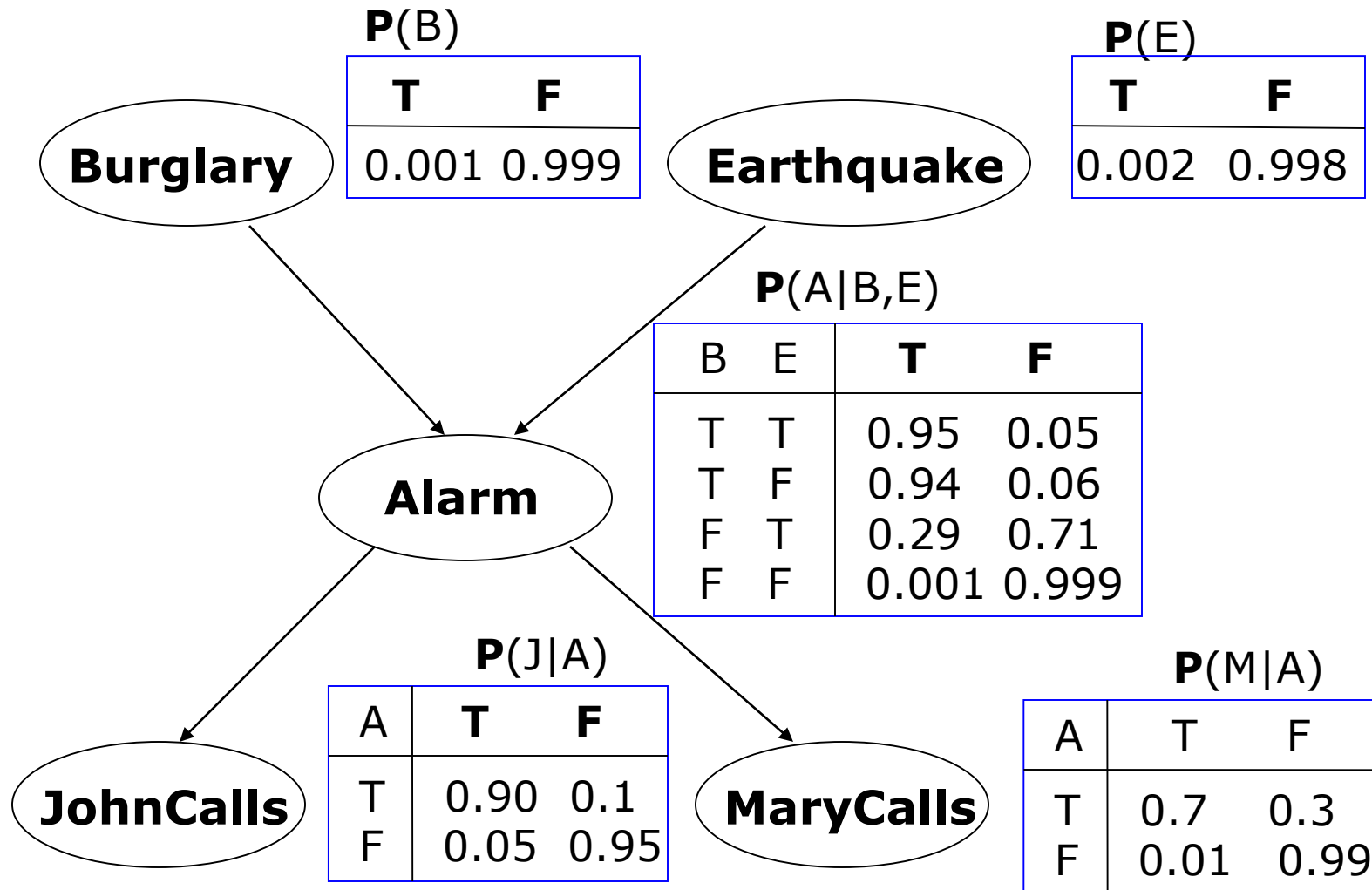
# Bayesian belief network

1. **Directed acyclic graph**: **Nodes** = random variables  
**Links** = dependencies between variables (missing links imply conditional independences among variables)
2. **Local conditional distributions**



# Bayesian belief network

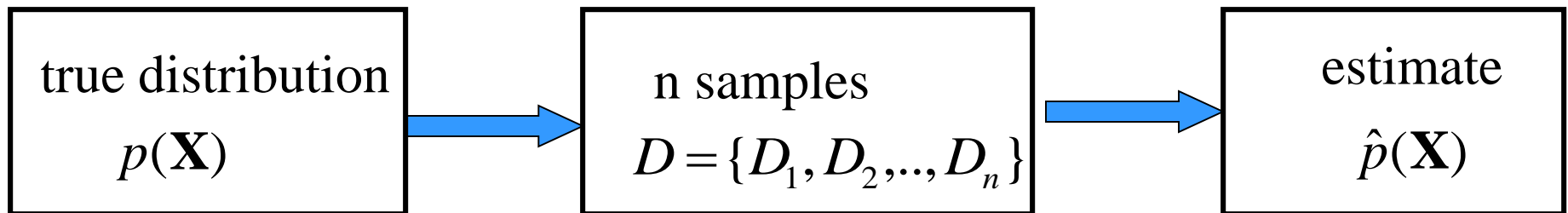
- How to obtain the parameters of the BBNs?



# Density estimation

**Data:**  $D = \{D_1, D_2, \dots, D_n\}$   
 $D_i = \mathbf{x}_i$  a vector of attribute values

**Objective:** try to estimate the underlying true probability distribution over variables  $\mathbf{X}$ ,  $p(\mathbf{X})$ , using examples in  $D$



**Standard (iid) assumptions:** Samples

- are **independent** of each other
- come from the same **(identical) distribution** (fixed  $p(\mathbf{X})$ )

# Learning via parameter estimation

In this lecture we consider **parametric density estimation**

## Basic settings:

- A set of random variables  $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$
- **A model of the distribution** over variables in  $\mathbf{X}$   
with parameters  $\Theta$
- **Data**  $D = \{D_1, D_2, \dots, D_n\}$

**Objective:** find parameters  $\hat{\Theta}$  that fit the data the best

For BBNs the parameters are local conditional probabilities

- What is the best set of parameters?
  - There are various criteria one can apply here.

# Parameter estimation. Basic criteria.

- **Maximum likelihood (ML)**

$$\text{maximize } p(D | \Theta, \xi)$$

$\xi$  - represents prior (background) knowledge

- **Maximum a posteriori probability (MAP)**

$$\text{maximize } p(\Theta | D, \xi)$$

**Selects the mode of the posterior**

$$p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{p(D | \xi)}$$

# Parameter estimation. Coin example.

**Coin example:** we have a coin that can be biased

**Outcomes:** two possible values -- head or tail

**Data:**  $D$  a sequence of outcomes  $x_i$  such that

- **head**  $x_i = 1$
- **tail**  $x_i = 0$

**Model:** probability of a head  $\theta$   
probability of a tail  $(1 - \theta)$

**Objective:**

We would like to estimate the probability of a **head**  $\hat{\theta}$   
from data

# Parameter estimation. Example.

- **Assume** the unknown and possibly biased coin
- Probability of the head is  $\theta$

- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

– **Heads:** 15

– **Tails:** 10

What would be your estimate of the probability of a head ?

$$\tilde{\theta} = ?$$



# Parameter estimation. Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is  $\theta$

- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

– **Heads:** 15

– **Tails:** 10

What would be your choice of the probability of a head ?

**Solution:** use frequencies of occurrences to do the estimate

$$\tilde{\theta} = \frac{15}{25} = 0.6$$

This is **the maximum likelihood estimate** of the parameter  $\theta$

# Probability of an outcome

**Data:**  $D$  a sequence of outcomes  $x_i$  such that

- **head**  $x_i = 1$
- **tail**  $x_i = 0$

**Model:** probability of a head  $\theta$   
probability of a tail  $(1 - \theta)$

**Assume:** we know the probability  $\theta$

**Probability of an outcome of a coin flip**  $x_i$

$$P(x_i | \theta) = \theta^{x_i} (1 - \theta)^{(1-x_i)} \quad \leftarrow \quad \text{Bernoulli distribution}$$

- Combines the probability of a head and a tail
- So that  $x_i$  is going to pick its correct probability
- Gives  $\theta$  for  $x_i = 1$
- Gives  $(1 - \theta)$  for  $x_i = 0$

# Probability of a sequence of outcomes.

**Data:**  $D$  a sequence of outcomes  $x_i$  such that

- **head**  $x_i = 1$
- **tail**  $x_i = 0$

**Model:** probability of a head  $\theta$   
probability of a tail  $(1 - \theta)$

**Assume:** a sequence of independent coin flips

**D = H H T H T H** (encoded as **D = 110101**)

What is the probability of observing the data sequence **D**:

$$P(D | \theta) = ?$$

# Probability of a sequence of outcomes.

**Data:**  $D$  a sequence of outcomes  $x_i$  such that

- **head**  $x_i = 1$
- **tail**  $x_i = 0$

**Model:** probability of a head  $\theta$   
probability of a tail  $(1 - \theta)$

**Assume:** a sequence of coin flips  $D = \text{H H T H T H}$   
encoded as  $D = 110101$

What is the probability of observing a data sequence  $D$ :

$$P(D \mid \theta) = \theta\theta(1 - \theta)\theta(1 - \theta)\theta$$

# Probability of a sequence of outcomes.

**Data:**  $D$  a sequence of outcomes  $x_i$  such that

- **head**  $x_i = 1$
- **tail**  $x_i = 0$

**Model:** probability of a head  $\theta$   
probability of a tail  $(1 - \theta)$

**Assume:** a sequence of coin flips  $D = \text{H H T H T H}$   
encoded as  $D = 110101$

What is the probability of observing a data sequence  $D$ :

$$P(D | \theta) = \theta\theta(1 - \theta)\theta(1 - \theta)\theta$$

 **likelihood of the data**

# Probability of a sequence of outcomes.

**Data:**  $D$  a sequence of outcomes  $x_i$  such that

- **head**  $x_i = 1$
- **tail**  $x_i = 0$

**Model:** probability of a head  $\theta$   
probability of a tail  $(1 - \theta)$

**Assume:** a sequence of coin flips  $D = \text{H H T H T H}$   
encoded as  $D = 110101$

What is the probability of observing a data sequence  $D$ :

$$P(D \mid \theta) = \theta \theta (1 - \theta) \theta (1 - \theta) \theta$$

$$P(D \mid \theta) = \prod_{i=1}^6 \theta^{x_i} (1 - \theta)^{(1-x_i)}$$

Can be rewritten using the Bernoulli distribution:

# The goodness of fit to the data.

**Learning:** we do not know the value of the parameter  $\theta$

**Our learning goal:**

- Find the parameter  $\theta$  that fits the data  $D$  the best?

**One solution to the “best”:** Maximize the likelihood

$$P(D | \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)}$$

**Intuition:**

- more likely are the data given the model, the better is the fit

**Note:** Instead of an error function that measures how bad the data fit the model we have a measure that tells us how well the data fit :

$$Error(D, \theta) = -P(D | \theta)$$

# Maximum likelihood (ML) estimate.

**Likelihood of data:**

$$P(D \mid \theta, \xi) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)}$$

**Maximum likelihood** estimate

$$\theta_{ML} = \arg \max_{\theta} P(D \mid \theta, \xi)$$

**Optimize log-likelihood (the same as maximizing likelihood)**

$$l(D, \theta) = \log P(D \mid \theta, \xi) = \log \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)} =$$
$$\sum_{i=1}^n x_i \log \theta + (1 - x_i) \log(1 - \theta) = \log \theta \underbrace{\sum_{i=1}^n x_i}_{N_1} + \log(1 - \theta) \underbrace{\sum_{i=1}^n (1 - x_i)}_{N_2}$$

$N_1$  - number of heads seen       $N_2$  - number of tails seen



# Maximum likelihood (ML) estimate.

## Optimize log-likelihood

$$l(D, \theta) = N_1 \log \theta + N_2 \log(1 - \theta)$$

## Set derivative to zero

$$\frac{\partial l(D, \theta)}{\partial \theta} = \frac{N_1}{\theta} - \frac{N_2}{(1 - \theta)} = 0$$

## Solving

$$\theta = \frac{N_1}{N_1 + N_2}$$

<b>ML Solution:</b> $\theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$
---------------------------------------------------------------------------

# Maximum likelihood estimate. Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is  $\theta$

- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

– **Heads:** 15

– **Tails:** 10

What is the ML estimate of the probability of a head and a tail?

# Maximum likelihood estimate. Example

- Assume the unknown and possibly biased coin
- Probability of the head is  $\theta$
- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

– **Heads:** 15

– **Tails:** 10

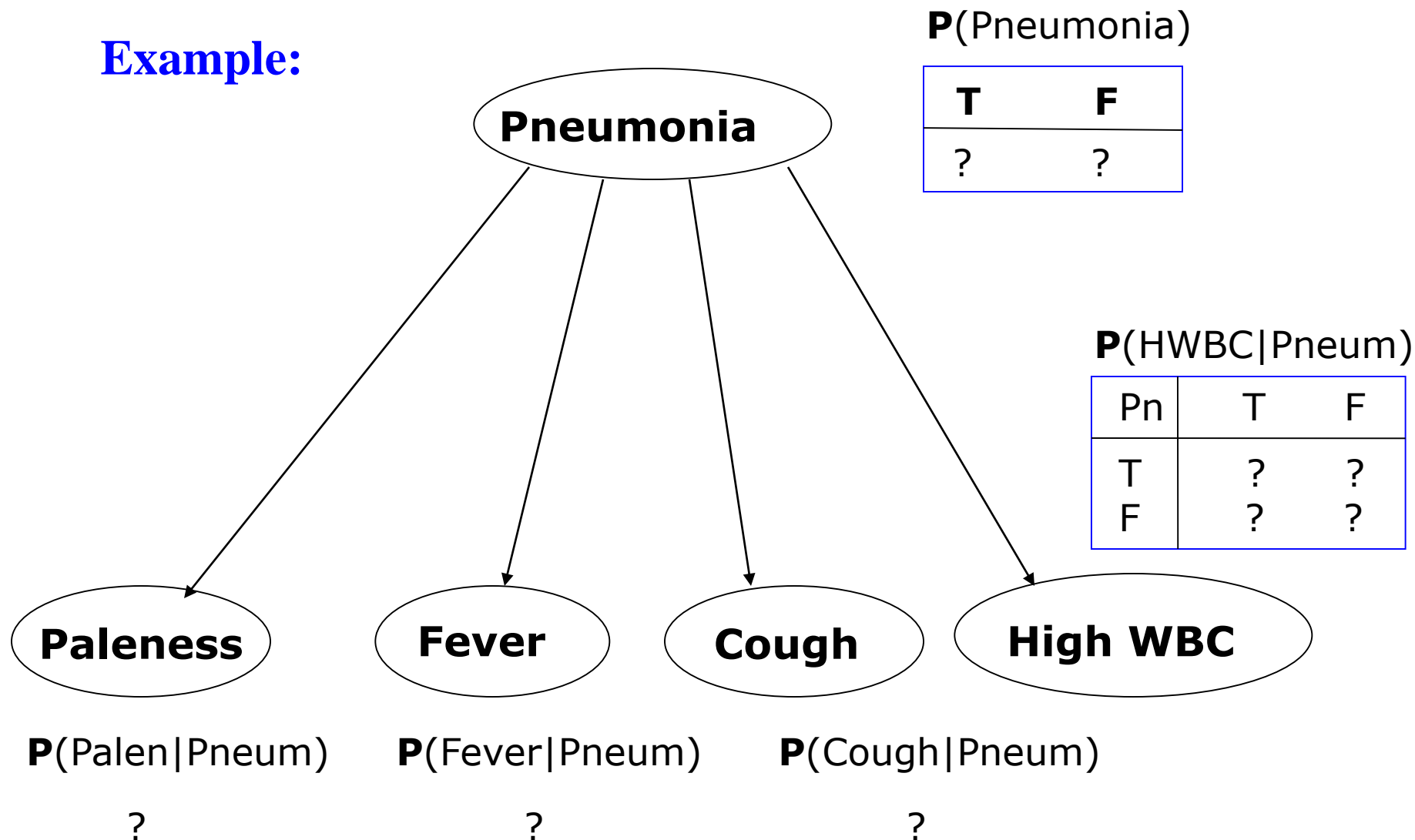
What is the ML estimate of the probability of head and tail ?

**Head:** 
$$\theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} = \frac{15}{25} = 0.6$$

**Tail:** 
$$(1 - \theta_{ML}) = \frac{N_2}{N} = \frac{N_2}{N_1 + N_2} = \frac{10}{25} = 0.4$$

# Learning of BBN parameters. Example.

Example:

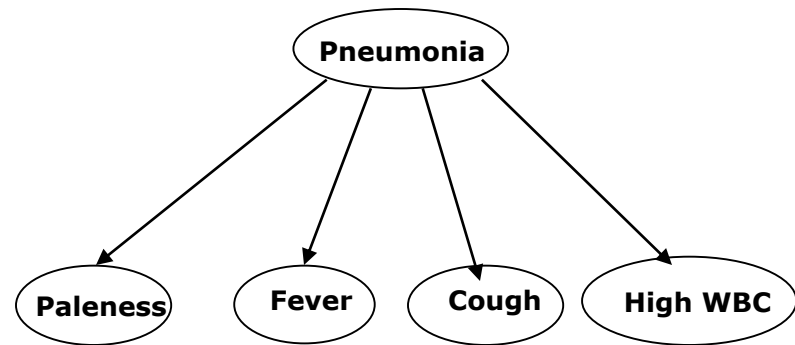


# Learning of BBN parameters. Example.

**Data D (different patient cases):**

**Pal Fev Cou HWB Pneu**

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F



# Estimates of parameters of BBN

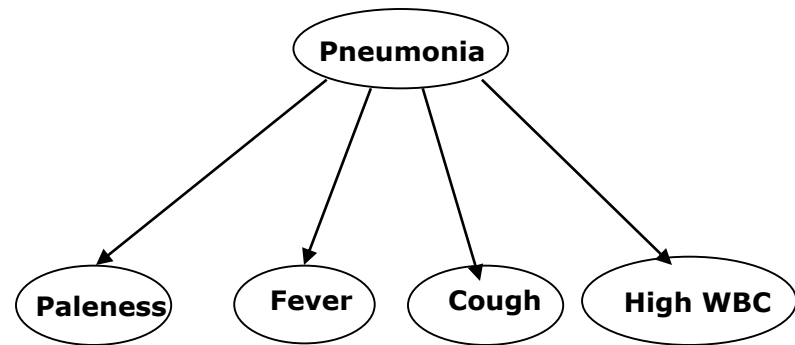
- Much like multiple **coin tosses**
- A “smaller” learning problem corresponds to the learning of exactly one conditional distribution
- **Example:**  
$$\mathbf{P}(\textit{Fever} \mid \textit{Pneumonia} = T)$$
- **Problem:** How to pick the data to learn?

# Learning of BBN parameters. Example.

**Data D (different patient cases):**

**Pal Fev Cou HWB Pneu**

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F



**How to estimate:**

$$\mathbf{P}(Fever | Pneumonia = T) = ?$$

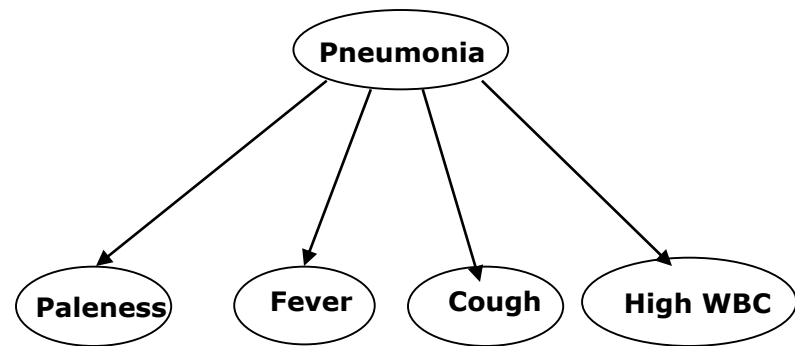
# Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 1:** Select data points with Pneumonia=T

Pal Fev Cou HWB Pneu

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F





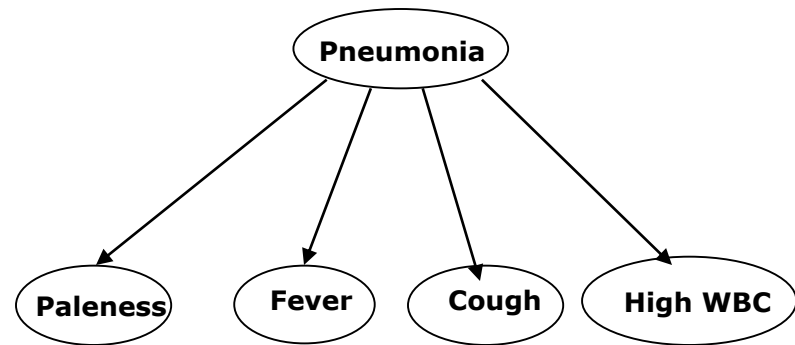
# Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 1:** Ignore the rest

**Pal Fev Cou HWB Pneu**

<b>F</b>	<b>F</b>	<b>T</b>	<b>T</b>	<b>T</b>
<b>F</b>	<b>F</b>	<b>T</b>	<b>F</b>	<b>T</b>
<b>F</b>	<b>T</b>	<b>T</b>	<b>T</b>	<b>T</b>
<b>T</b>	<b>T</b>	<b>T</b>	<b>T</b>	<b>T</b>
<b>F</b>	<b>T</b>	<b>F</b>	<b>T</b>	<b>T</b>



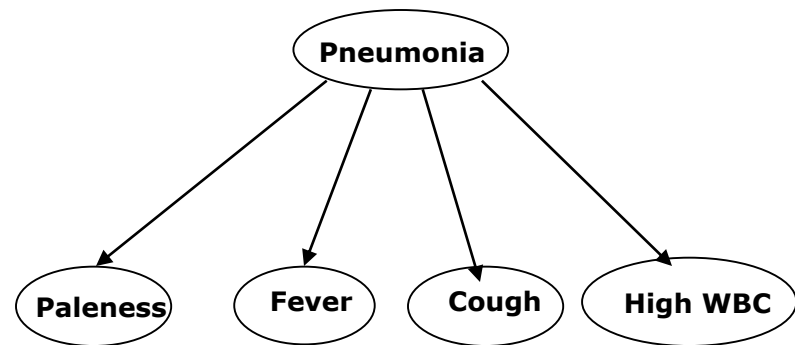
# Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 2:** Select values of the random variable defining the distribution of Fever

**Pal** **Fev** **Cou** **HWB** **Pneu**

F	<b>F</b>	T	T	T
F	<b>F</b>	T	F	T
F	<b>T</b>	T	T	T
T	<b>T</b>	T	T	T
F	<b>T</b>	F	T	T



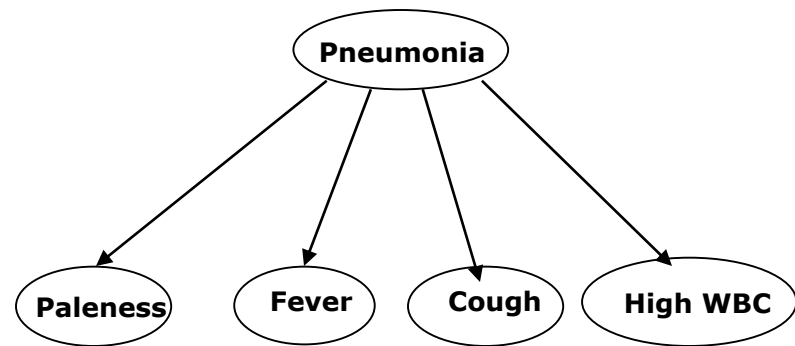
# Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 2:** Ignore the rest

**Fev**

**F**  
**F**  
**T**  
**T**  
**T**



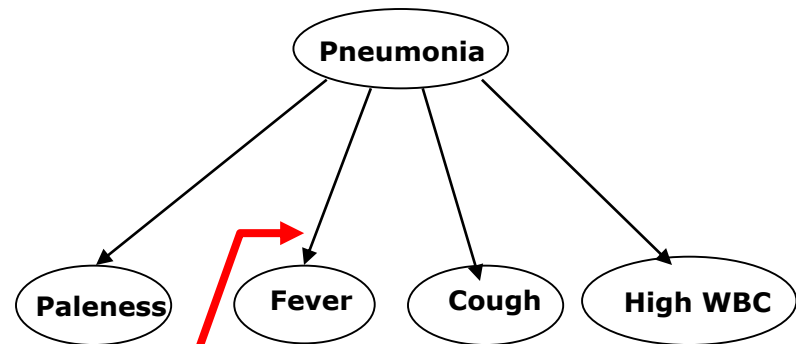
# Learning of BBN parameters. Example.

**Learn:**  $P(\text{Fever} | \text{Pneumonia} = T)$

**Step 3: Learning the ML estimate**

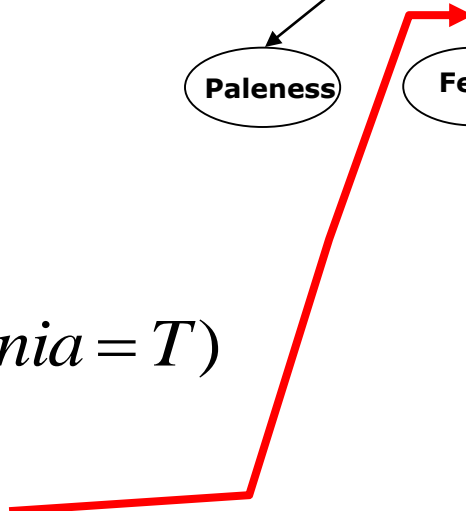
**Fev**

**F**  
**F**  
**T**  
**T**  
**T**



$P(\text{Fever} | \text{Pneumonia} = T)$

<b>T</b>	<b>F</b>
0.6	0.4



# Types of learning

- **Supervised learning**
  - Learning mapping between input  $\mathbf{x}$  and desired output  $\mathbf{y}$
  - Teacher gives me  $\mathbf{y}$ 's for the learning purposes
- **Unsupervised learning**
  - Learning relations between data components
  - No specific outputs given by a teacher
- **Reinforcement learning**
  - Learning mapping between input  $\mathbf{x}$  and desired output  $\mathbf{y}$
  - Critic does not give me  $\mathbf{y}$ 's but instead a signal (reinforcement) of how good my answer was
- **Other types of learning:**
  - **Concept learning, Active learning, Transfer learning, Deep learning**

# Supervised learning

**Data:**  $D = \{D_1, D_2, \dots, D_n\}$  a set of  $n$  examples

$$D_i = \langle \mathbf{x}_i, y_i \rangle$$

$\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,d})$  is an input vector of size  $d$

$y_i$  is the desired output (given by a teacher)

**Objective:** learn the mapping  $f : X \rightarrow Y$

s.t.  $y_i \approx f(\mathbf{x}_i)$  for all  $i = 1, \dots, n$

- **Regression:**  $Y$  is **continuous**

Example: earnings, product orders  $\rightarrow$  company stock price

- **Classification:**  $Y$  is **discrete**

Example: handwritten digit in binary form  $\rightarrow$  digit label