# MONASH University
## Information Technology

# FIT 3080: Intelligent Systems

## Probability

## Chapter 13

Many slides are adapted from Stuart Russell,  Andrew Moore
or Dan Klein

# So far

**Agents did not consider:**

- **Uncertainty about the world or the outcome of an action**
- **Learning their knowledge**

MONASH University
Information Technology

# From Now On

- **Uncertainty**
  - Probability, Bayesian Networks
- **Planning for Complex Decisions**
  - Markov Decision Processes, Reinforcement Learning
- **Machine Learning**
  - Classification, Regression

MONASH University
Information Technology

# Outline

- **Background:**
  - Random variables and probabilistic inference
  - Probabilistic models
  - Joint, marginal and conditional distributions
- **Inference by enumeration**
- **Product Rule, Chain Rule, Bayes' Rule**
- **Independence and conditional independence**

MONASH University
Information Technology
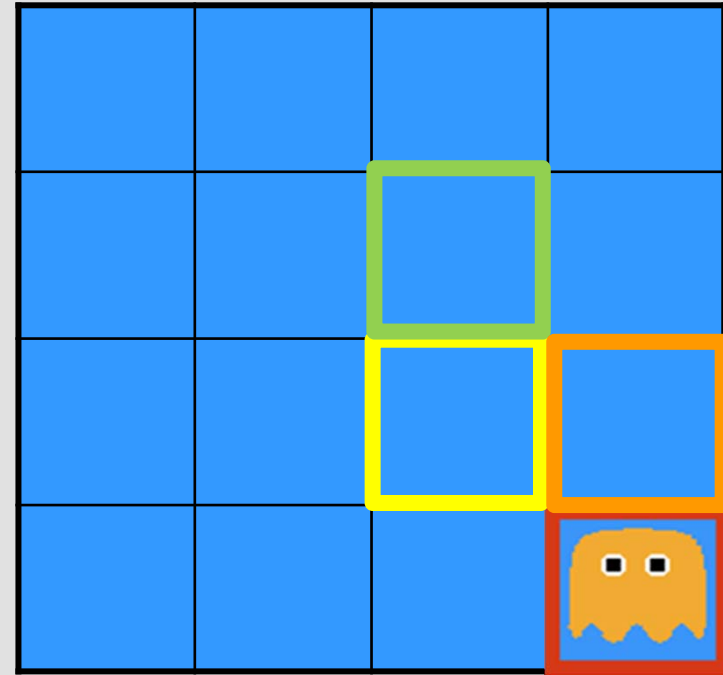
# Random Variables

- **A random variable is some aspect of the world about which we may have uncertainty**
  - R = Is it raining?
  - D = How long will it take to drive to work?
  - L = Where am I?
- **We denote random variables with capital letters**
- **Random variables have domains**
  - R in {true, false}   (sometimes write as {+r, $\neg$r})
  - D in [0, $\infty$)
  - L in possible locations, maybe {(0,0), (0,1), …}

# Probabilistic Inference

- **Probabilistic inference: compute a desired probability from other known probabilities**
- **We generally compute conditional probabilities**
  - They represent an agent's *beliefs* given the evidence
  - E.g., Pr(on time | no reported accidents) = 0.90
- **Probabilities change with new evidence:**
  - Observing new evidence causes *beliefs to be updated*
  - E.g., Pr(on time | no accidents, 5 a.m.) = 0.95
    
    Pr(on time | no accidents, 5 a.m., raining) = 0.80

MONASH University
Information Technology

# Example – Inference in Ghostbusters

- **A ghost is somewhere in the grid**
- **Sensor readings tell how close a tile is to the ghost**
  - On the ghost: **red**
  - 1 away: **orange**
  - 2 away: **yellow**
  - 3+ away: **green**
- **Sensors are noisy, but we know Pr(Color|Distance)**



| Pr(red\|2) | Pr(orange\|2) | Pr(yellow\|2) | Pr(green\|2) |
|------------|---------------|---------------|--------------|
| 0.05 | 0.17 | 0.46 | 0.32 |

**We want to know: Pr(Location | Color)**

MONASH University
Information Technology

# Uncertainty and Probabilistic Inference

- **General situation:**
  - **Evidence**: Agent knows certain things about the state of the world
  - **Hidden variables**: Agent needs to reason about other aspects
  - **Model**: Agent knows something about how the known variables relate to the unknown variables
- **Probabilistic reasoning gives us a framework for managing our beliefs and knowledge**

No observations

| | | |
|---|---|---|
| 0.11 | 0.11 | 0.11 |
| 0.11 | 0.11 | 0.11 |
| 0.11 | 0.11 | 0.11 |

Evidence: yellow

| | | |
|---|---|---|
| 0.17 | 0.10 | 0.10 |
| 0.09 | 0.17 | 0.10 |
| <0.01 | 0.09 | 0.17 |

Evidence: red

| | | |
|---|---|---|
| <0.01 | <0.01 | 0.03 |
| <0.01 | 0.05 | 0.05 |
| <0.01 | 0.05 | 0.81 |

MONASH University
Information Technology

# Probabilistic Models (I)

- **Probabilistic models describe how (a portion of) the world works**
- **Models are always simplifications**
  - May not account for every variable
  - May not account for all interactions between variables
  - "All models are wrong; but some are useful."
    - George E. P. Box
- **What do we do with probabilistic models?**
  - We (or our agents) need to reason about unknown variables given evidence
    - > explanation (diagnostic reasoning)
    - > prediction (causal reasoning)
    - > value of information

# Probability Distributions

- **Unobserved random variables have distributions that represent probabilities of value assignments**

Pr(Temp)

| Temp | Pr |
|------|-----|
| warm | 0.5 |
| cold | 0.5 |

Pr(Weather)

| Weather | Pr |
|---------|-----|
| sunny | 0.6 |
| rain | 0.1 |
| fog | 0.3 |

- **A probability is a single number**

$$\Pr(\text{Weather}=\text{rain}) = 0.1 \quad \text{or} \quad \Pr(\text{rain}) = 0.1$$

- **Kolmogorov's axioms:**

$$\forall x \quad \Pr(x) \geq 0 \qquad \sum_x \Pr(x) = 1$$

MONASH University
Information Technology

# Joint Distributions

- **A *joint distribution* over a set of random variables $X_1, \ldots, X_n$ specifies a real number for each value assignment (or *outcome*):**

  $$\Pr(X_1{=}x_1, \ldots, X_n{=}x_n) \text{ or } \Pr(x_1, \ldots, x_n)$$

  - Size of distribution of n variables with domain sizes d?

- **Must obey:**

  $$\forall x_i \quad \Pr(x_1, \ldots, x_n) \geq 0$$

  $$\sum_{x_1, \ldots, x_n} \Pr(x_1, \ldots, x_n) = 1$$

- **For all but small distributions, impractical to write out**

$$\Pr(W, T)$$

| T | W | Pr |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

MONASH University
Information Technology
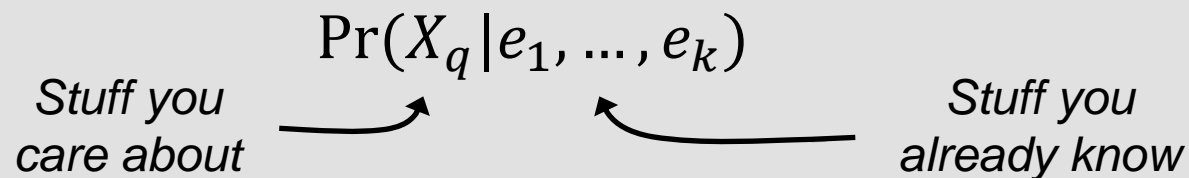
# Probabilistic Models (II)

- **A probabilistic model is a joint distribution over a set of variables**

$$\Pr(X_1, X_2, \ldots, X_n)$$

- **Given a joint distribution, we can reason about unobserved variables given evidence**

- **General form of a query:**

$$\Pr(X_q | e_1, \ldots, e_k)$$

*Stuff you care about*

*Stuff you already know*

- **This kind of _posterior distribution_ is also called the _belief function_ of an agent who uses this model**

MONASH University
Information Technology

# Events

- **An _outcome_ is a possible result of an experiment**
- **An _event_ is a set _E_ of outcomes**

$$\Pr(E) = \sum_{\{x_1,\ldots,x_n\} \in E} \Pr(x_1, \ldots, x_n)$$

- **From a joint distribution, we can calculate the probability of any event**
  - Probability that it is hot AND sunny
  - Probability that it is hot
  - Probability that it is hot OR sunny
- **Typically, the events we care about are _partial assignments_, like Pr(T=hot)**

$\mathbf{Pr}(\boldsymbol{W}, \boldsymbol{T})$

| T | W | Pr |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

MONASH University
Information Technology

# Marginal Distributions

- **_Marginal distributions_ are sub-tables that eliminate variables**
- **_Marginalization_ (summing out): Combine collapsed rows by adding**

$\mathbf{Pr}(W, T)$

| T | W | Pr |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

$$\Pr(t) = \sum_{w \in \{sun, rain\}} \Pr(t, w)$$

$$\Pr(w) = \sum_{t \in \{hot, cold\}} \Pr(t, w)$$

$\Pr(T)$

| T | Pr |
|------|-----|
| hot | 0.5 |
| cold | 0.5 |

$\Pr(W)$

| W | Pr |
|------|-----|
| sun | 0.6 |
| rain | 0.4 |

# Conditional Distributions (I)

- **Conditional distributions are probability distributions over some variables given fixed values of others**

**Joint Distribution**

$$\mathbf{Pr}(W, T)$$

| T | W | Pr |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

**Conditional Distributions**

$$\mathbf{Pr}(W|T = hot)$$

| W | Pr |
|------|-----|
| sun | 0.8 |
| rain | 0.2 |

$$\mathbf{Pr}(W|T = cold)$$

| W | Pr |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

$\mathbf{Pr}(W|T)$

MONASH University
Information Technology

# Conditional Distributions (II)

$$\Pr(X \mid Y) = \frac{\Pr(X \wedge Y)}{\Pr(Y)}$$

$$\Pr(X \cap Y)$$



$$\Pr(Y) \qquad \Pr(X)$$

**Pr(*W*, *T*)**

| T | W | Pr |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

$$\Pr(W = rain | T = cold) = ?$$

MONASH University
Information Technology

# Conditional Distributions (III)

- ***Conditional* or *posterior probabilities:***
  - E.g., Pr(*cavity* | *toothache*)=0.8, given that *toothache* is all I know
- **Notation for conditional distributions:**
  - Pr(*cavity* | *toothache*) = a single number
  - Pr(Cavity, Toothache) = 2x2 table sums to 1
  - Pr(Cavity | Toothache) = Two 2-element vectors, each sums to 1
- **If we know more:**
  - Pr(*cavity* | *toothache*, *catch*) = 0.9
  - Pr(*cavity* | *toothache*, *cavity*) = 1
- **Less specific beliefs remain *valid* after more evidence arrives, but are not always *useful***
- **New evidence may be irrelevant, allowing simplification:**
  - Pr(*cavity* | *toothache*, *traffic*) = P(*cavity* | *toothache*) = 0.8

MONASH University
Information Technology

# Normalization Trick

- **A trick to get a whole conditional distribution at once:**
  - Select the joint probabilities matching the evidence
  - *Normalize* the selection (make it sum to one)

$\mathbf{Pr}(W, T)$

| T | W | Pr |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

→ Select →

$\mathrm{Pr}(T, rain)$

| T | R | Pr |
|------|------|-----|
| hot | rain | 0.1 |
| cold | rain | 0.3 |

→ Normalize →

$\mathrm{Pr}(T \mid rain)$

| T | Pr |
|------|------|
| hot | 0.25 |
| cold | 0.75 |

- **Why does this work?**

$$\mathrm{Pr}(x_1 \mid x_2) = \frac{\mathrm{Pr}(x_1, x_2)}{\mathrm{Pr}(x_2)} = \frac{\mathrm{Pr}(x_1, x_2)}{\sum_{x_1} \mathrm{Pr}(x_1, x_2)}$$

# Inference by Enumeration (I)

- **Pr(sun)?**

- **Pr(sun | summer)?**

- **Pr(sun | winter, hot)?**

| S | T | W | Pr |
|--------|------|------|------|
| summer | hot | sun | 0.30 |
| summer | hot | rain | 0.05 |
| summer | cold | sun | 0.10 |
| summer | cold | rain | 0.05 |
| winter | hot | sun | 0.10 |
| winter | hot | rain | 0.05 |
| winter | cold | sun | 0.15 |
| winter | cold | rain | 0.20 |

# Inference by Enumeration (II)

- **General case:**
  - Evidence variables:   $E_1, \ldots, E_k = e_1, \ldots, e_k$
  - Query variable(s):   $Q$ $\left. \vphantom{\begin{array}{c}a\\b\\c\end{array}} \right\}$ $X_1, \ldots, X_n$
  - Unknown variables:   $U_1, \ldots, U_r$ *All variables*
- **We want** $\Pr(\boldsymbol{Q}|\boldsymbol{e_1}, \ldots, \boldsymbol{e_k})$
- **Procedure**
  1. Select the entries that are consistent with the evidence
  2. Sum out *U* to get the joint probability of Query and Evidence:
  $$\Pr(Q, e_1, \ldots, e_k) = \sum_{u_1, \ldots, u_r} \Pr(\underbrace{Q, u, \ldots, u_r, e_1, \ldots, e_k}_{X_1, \ldots, X_n})$$

  3. Normalize the remaining entries to conditionalize
- **Problems:**
  - Worst-case time complexity $O(d^n)$
  - Space complexity $O(d^n)$ to store the joint distribution

# Inference by Enumeration – Example

- **Pr(sun | summer)**
  - Evidence variables?
  - Query variables?
  - Unknown variables?

- **Procedure**
  - Select entries
  - Sum out *U* to get a joint probability of *Q* and *E*
  - Normalize the remaining entries to conditionalize

| S | T | W | Pr |
|---|---|---|---|
| summer | hot | sun | 0.30 |
| summer | hot | rain | 0.05 |
| summer | cold | sun | 0.10 |
| summer | cold | rain | 0.05 |
| winter | hot | sun | 0.10 |
| winter | hot | rain | 0.05 |
| winter | cold | sun | 0.15 |
| winter | cold | rain | 0.20 |

# The Product Rule

- **Sometimes we have conditional distributions but want the joint distribution**

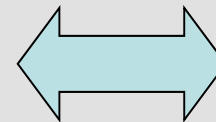$$\Pr(x|y) = \frac{\Pr(x,y)}{\Pr(y)} \quad \Longleftrightarrow \quad \Pr(x,y) = \Pr(x|y)\Pr(y)$$

- **Example:**

**Pr(W)**

| W | Pr |
|------|-----|
| sun | 0.8 |
| rain | 0.2 |

**Pr(D|W)**

| D | W | Pr |
|-----|------|-----|
| wet | sun | 0.1 |
| dry | sun | 0.9 |
| wet | rain | 0.7 |
| dry | rain | 0.3 |

**Pr(D, W)**

| D | W | Pr |
|-----|------|------|
| wet | sun | 0.08 |
| dry | sun | 0.72 |
| wet | rain | 0.14 |
| dry | rain | 0.06 |

MONASH University
Information Technology

# The Chain Rule

- **We can always write a joint distribution as an incremental product of conditional distributions**

$$\Pr(x_1, \ldots, x_n) = \prod_{i=1}^{n} \Pr(x_i | x_1, \ldots, x_{i-1})$$

- **Example:**
Pr(Traffic,Umbrella,Rain)=
    Pr(Umbrella|Rain,Traffic) x Pr(Traffic|Rain) x Pr(Rain)

- **Why is this true?**

# Bayes' Rule

- **Two ways to factor a joint distribution over two variables:**

$$\Pr(x, y) = \Pr(x|y)\Pr(y) = \Pr(y|x)\Pr(x)$$

$$\Pr(x|y) = \frac{\Pr(y|x)\Pr(x)}{\Pr(y)}$$

- **Why is this helpful?**
  – Lets us build one conditional from its reverse
  – Often one conditional is tricky but the other one is simple
  – Foundation of many systems (e.g., ASR, MT)

# Bayes Rule: Conditionalization

- **Attributed to Rev. Thomas Bayes**

$$\Pr(h \mid e) = \frac{\Pr(e \mid h) \Pr(h)}{\Pr(e)}$$

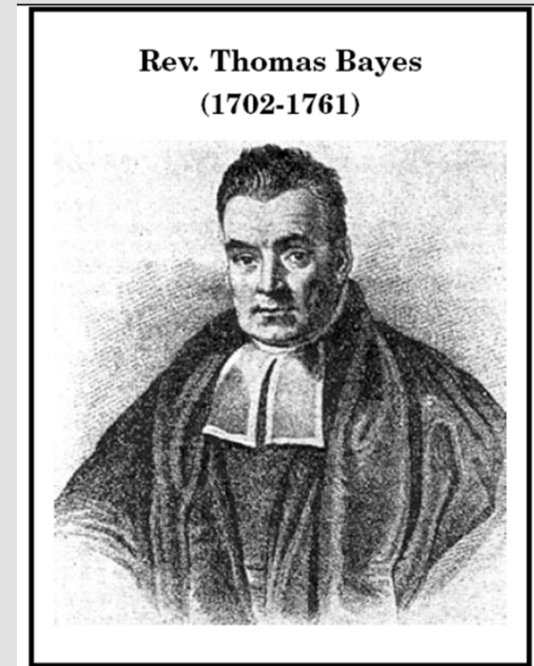- **Also called *Conditionalization*:**

$$\Pr'(h) = \Pr(h \mid e)$$

- **Also read as**

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Prob of evidence}}$$

- **Assumptions:**
  - Joint priors over $\{h_i\}$ and $e$ exist
  - Total evidence: $e$ is observed

Rev. Thomas Bayes
(1702-1761)

# Inference with Bayes' Rule – Example

**Diagnosis of breast cancer (hypothesis), given xray (evidence)**

- **Let $\Pr(h){=}0.01$, $\Pr(e/h){=}0.8$ and $\Pr(e/{\sim}h){=}0.1$**
- **Bayes theorem yields**

$$\Pr(h \mid e) = \frac{\Pr(e \mid h)\Pr(h)}{\Pr(e)}$$

$$= \frac{\Pr(e \mid h)\Pr(h)}{\Pr(e \mid h)\Pr(h) + \Pr(e \mid {\sim} h)\Pr({\sim} h)}$$

$$= \frac{0.8 \times 0.01}{0.8 \times 0.01 + 0.1 \times 0.99}$$

$$= \frac{0.008}{0.008 + 0.099} = \frac{0.008}{0.107} \approx 0.075$$

MONASH University
Information Technology

# Ghostbusters Revisited

- **We have two distributions:**
  - **Prior distribution** over ghost location: $\Pr(L)$
  - **Sensor model**: $\Pr(R \mid D)$
    - > Given by some "black box" process
    - > Assume reading is at the lower left corner
    - > E.g., Pr(yellow|D≥3)=0.27
      Pr(yellow|D=2)=0.46
      Pr(yellow|D=1)=0.25
      Pr(yellow|D=0)=0.03

| 0.11 | 0.11 | 0.11 |
|------|------|------|
| 0.11 | 0.11 | 0.11 |
| 0.11 | 0.11 | 0.11 |

- **The posterior distribution Pr(L|R) over ghost locations given a reading**

$$\Pr(l = (3,1)|yellow)$$
$$\propto \Pr(yellow|l = (3,1))\Pr(l = (3,1))$$
$$\propto 0.03 * 0.11 = 0.0033$$

| 0.17 | 0.10 | 0.10 |
|------|------|------|
| 0.09 | 0.17 | 0.10 |
| <0.01 | 0.09 | 0.17 |

MONASH University
Information Technology

# Example Problems

- Suppose a murder occurs in a town of population 10,000 (10,001 before the murder). A suspect is brought in and DNA tested. The probability that there is a DNA match give that a person is innocent is 1/100,000; the probability of a match on a guilty person is 1. What is the probability he is guilty given a DNA match?

- Doctors have found that people with Creutzfeldt–Jakob disease (CJ) almost invariably ate lots of hamburgers, thus Pr(HamburgerEater|CJ) = 0.9. CJ is a rare disease: about 1 in 100,000 people get it. Eating hamburgers is widespread: Pr(HamburgerEater) = 0.5. What is the probability that a regular hamburger eater will have CJ disease?

MONASH University
Information Technology

# Independence

- **Two variables are _independent_ if:**
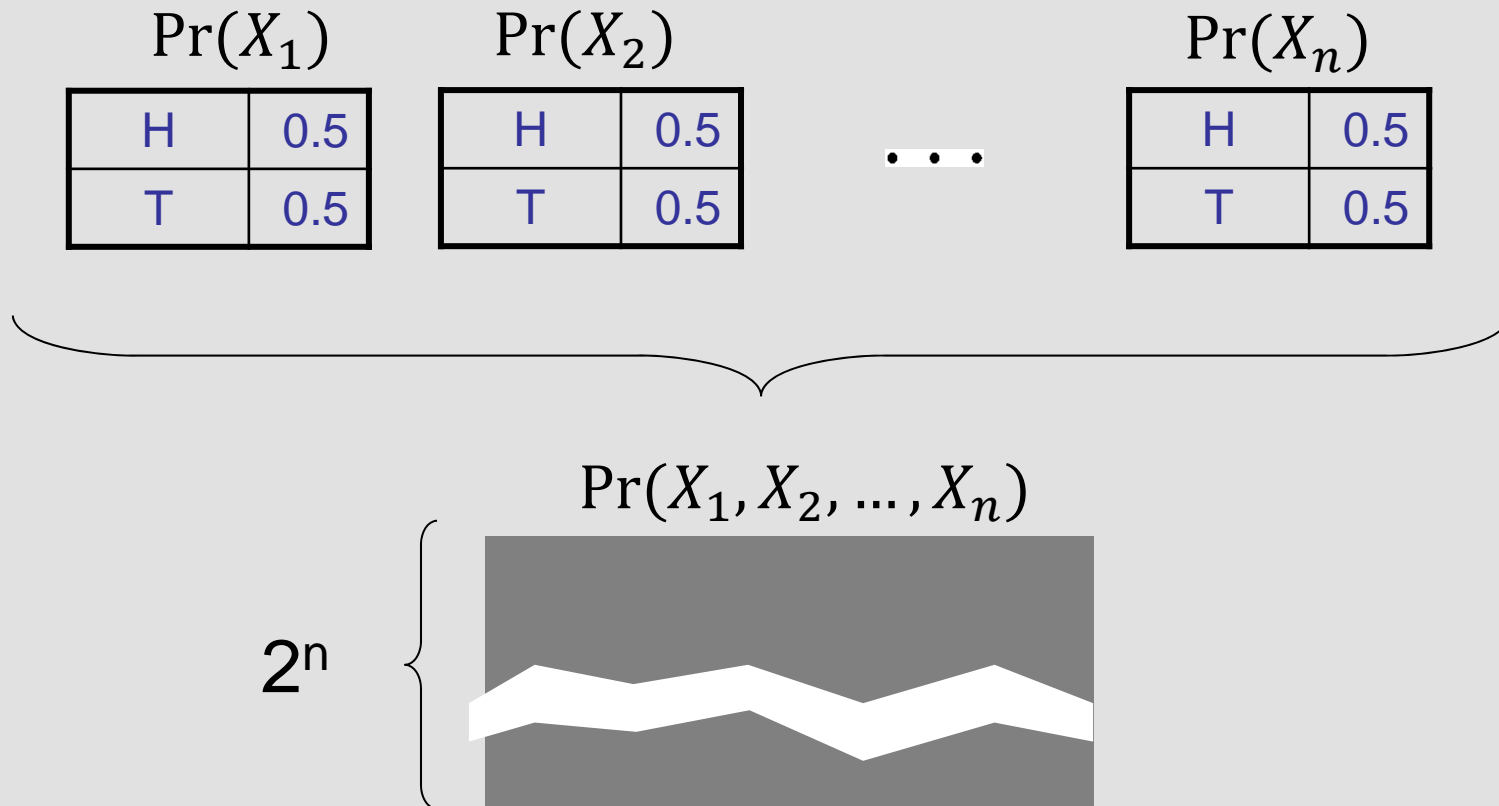
$$\Pr(X, Y) = \Pr(X)\Pr(Y)$$

$$\forall x, y \quad \Pr(x, y) = \Pr(x)\Pr(y) \quad \text{or} \quad \Pr(x|y) = \Pr(x)$$

$$X \perp\!\!\!\perp Y$$

- **Independence is a simplifying _modeling assumption_**
  - _Empirical_ joint distributions: at best "close" to independent
  - What could we assume for
    {Weather, Traffic, Cavity, Toothache}?

# Independence – Example

- **N fair, independent coin flips:**

$$\Pr(X_1) \qquad \Pr(X_2) \qquad\qquad\qquad \Pr(X_n)$$

| H | 0.5 |
|---|-----|
| T | 0.5 |

| H | 0.5 |
|---|-----|
| T | 0.5 |

. . .

| H | 0.5 |
|---|-----|
| T | 0.5 |

$$\Pr(X_1, X_2, \ldots, X_n)$$

$2^n$

# Which Variables are Independent?

**Pr(T)**

| T | Pr |
|------|-----|
| warm | 0.5 |
| cold | 0.5 |

**Pr₁(T, W)**

| T | W | Pr |
|------|------|-----|
| warm | sun | 0.4 |
| warm | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

**Pr₂(T, W)**

| T | W | Pr |
|------|------|-----|
| warm | sun | 0.3 |
| warm | rain | 0.2 |
| cold | sun | 0.3 |
| cold | rain | 0.2 |

**Pr(W)**

| W | Pr |
|------|-----|
| sun | 0.6 |
| rain | 0.4 |

MONASH University
Information Technology

# Conditional Independence (I)

- **Employs domain knowledge to simplify probabilistic models**
- **Example: Pr(Toothache,Cavity,Catch)**
  **If I have or don't have a cavity, the probability that the probe catches in the tooth doesn't depend on whether I have a toothache:**
  - Pr(+catch | +toothache, +cavity) = Pr(+catch | +cavity)
  - Pr(+catch | +toothache, ¬cavity) = Pr(+catch| ¬cavity)
  ➔ Catch is **_conditionally independent_** of Toothache given Cavity:
    > Pr(Catch | Toothache, Cavity) = Pr(Catch | Cavity) or
    > Pr(Toothache | Catch, Cavity) = Pr(Toothache | Cavity) or
    > Pr(Toothache, Catch | Cavity) =
    $$Pr(Toothache | Cavity) \times Pr(Catch | Cavity)$$

# Conditional Independence (II)

- **Unconditional (absolute) independence is rare**
- **Conditional independence is our most basic and robust form of knowledge about uncertain environments:**

$$\forall x, y, z \qquad \Pr(x, y|z) = \Pr(x|z)\Pr(y|z) \text{ or}$$
$$\Pr(x|y, z) = \Pr(x|z)$$
$$\Pr(X, Y|Z) = \Pr(X|Z)\Pr(Y|Z)$$
$$\Pr(X|Y, Z) = \Pr(X|Z)$$

$$X \perp\!\!\!\perp Y \mid Z$$

- **Example**
Pr(Traffic,Umbrella|Rain)= Pr(Umbrella|Rain) x Pr(Traffic|Rain)  or
Pr(Traffic|Umbrella,Rain)=Pr(Traffic|Rain)
- **Bayesian networks / graphical models help us express conditional independence assumptions**

MONASH University
Information Technology

# Reading

- **Russell, S. and Norvig, P. (2010), *Artificial Intelligence – A Modern Approach* (3nd ed), Prentice Hall**
  - Chapter 13