# CSE 573: Artificial Intelligence
## Winter 2019

## Uncertainty &
## Probabilistic Reasoning

### Hanna Hajishirzi

# 573 Outline

- We're done with Part I: Search and Planning!

- Part II: Probabilistic Reasoning
  - Diagnosis
  - Speech recognition
  - Tracking objects
  - Robot mapping
  - Genetics
  - Error correcting codes
  - … lots more!

# Outline

- **Probability review**
  - Random Variables and Events
  - Joint / Marginal / Conditional Distributions
  - Product Rule, Chain Rule, Bayes' Rule
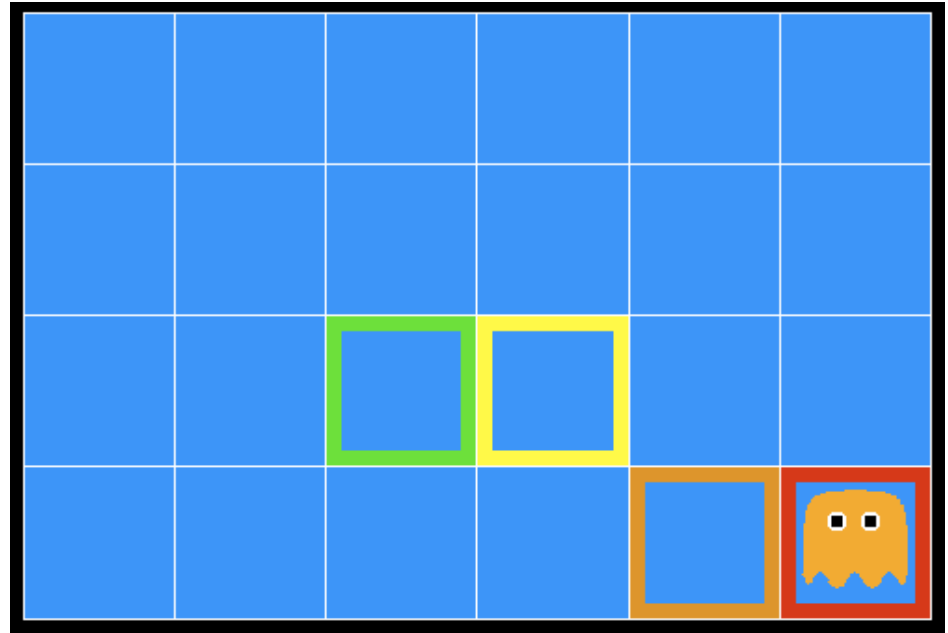  - Probabilistic Inference
  - Independence

# Probability Summary

- **Conditional probability**

$$P(x|y) = \frac{P(x,y)}{P(y)}$$

- **Product rule**

$$P(x,y) = P(x|y)P(y)$$

- **Chain rule**

$$P(X_1, X_2, \ldots X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1,X_2)\ldots$$
$$= \prod_{i=1}^{n} P(X_i|X_1,\ldots,X_{i-1})$$

- **X, Y independent if and only if:**   $\forall x, y : P(x,y) = P(x)P(y)$

- **X and Y are conditionally independent given Z if and only if:**   $X \perp\!\!\!\perp Y | Z$
$$\forall x, y, z : P(x,y|z) = P(x|z)P(y|z)$$

# Inference in Ghostbusters

- A ghost is in the grid somewhere

- Sensor readings tell how close a square is to the ghost
    - On the ghost: red
    - 1 or 2 away: orange
    - 3 or 4 away: yellow
    - 5+ away: green

- Sensors are noisy, but we know P(Color | Distance)

| P(red | 3) | P(orange | 3) | P(yellow | 3) | P(green | 3) |
|------------|----------------|----------------|----------------|
| 0.05 | 0.15 | 0.5 | 0.3 |

# 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}
  - D in [0, 1)
  - L in possible locations, maybe {(0,0), (0,1), …}

# Probability Distribution

- Unobserved random variables have distributions

$P(T)$

| T | P |
|------|-----|
| hot | 0.5 |
| cold | 0.5 |

$P(W)$

| W | P |
|--------|-----|
| sun | 0.6 |
| rain | 0.1 |
| fog | 0.3 |
| meteor | 0.0 |

Shorthand notation:

$$P(hot) = P(T = hot),$$
$$P(cold) = P(T = cold),$$
$$P(rain) = P(W = rain),$$
$$\ldots$$

OK *if* all domain entries are unique

- A distribution is a TABLE of probabilities of values

- A probability (lower case value) is a single number

$$P(W = rain) = 0.1$$

- Must have: $\forall x \ \ P(X = x) \geq 0$ and $\sum_x P(X = x) = 1$

# Joint Distributions

- A joint distribution over a set of random variables: $X_1, X_2, \ldots X_n$ specifies a real number for each *outcome* (ie each assignment):

$$P(X_1 = x_1, X_2 = x_2, \ldots X_n = x_n)$$

$$P(x_1, x_2, \ldots x_n)$$

  - Must obey:

$$P(x_1, x_2, \ldots x_n) \geq 0$$

$$\sum_{(x_1, x_2, \ldots x_n)} P(x_1, x_2, \ldots x_n) = 1$$

$P(T, W)$

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

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

- A probabilistic model is a joint distribution over variables of interest
- For all but the smallest distributions, impractical to write out

# Events

- An outcome is a joint assignment for all the variables

$$(x_1, x_2, \ldots x_n)$$

- An event is a set E of outcomes

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

- From a joint distribution, we can calculate the probability of any event

  - Probability that it's hot AND sunny?

  - Probability that it's hot?

  - Probability that it's hot OR sunny?

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

# Marginal Distributions

- Marginal distributions are sub-tables which eliminate variables
- Marginalization (summing out): Combine collapsed rows by adding

$$P(X_1 = x_1) = \sum_{x_2} P(X_1 = x_1, X_2 = x_2)$$

$P(T, W)$

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

$$P(t) = \sum_w P(t, w)$$

$$P(w) = \sum_t P(t, w)$$

$P(T)$

| T | P |
|------|-----|
| hot | 0.5 |
| cold | 0.5 |

$P(W)$

| W | P |
|------|-----|
| sun | 0.6 |
| rain | 0.4 |

# Quiz: Marginal Distribution

$P(X,Y)$

| X | Y | P |
|----|----|-----|
| +x | +y | 0.2 |
| +x | -y | 0.3 |
| -x | +y | 0.4 |
| -x | -y | 0.1 |

$$P(x) = \sum_y P(x,y)$$

$$P(y) = \sum_x P(x,y)$$

$P(X)$

| X | P |
|----|----|
| +x |  |
| -x |  |

$P(Y)$

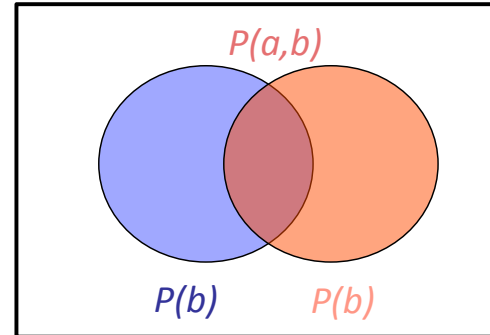| Y | P |
|----|----|
| +y |  |
| -y |  |

# Conditional Probability

- A simple relation between joint and conditional probabilities
  - In fact, this is taken as the *definition* of a conditional probability

$$P(a|b) = \frac{P(a,b)}{P(b)}$$



$P(a,b)$

$P(b)$     $P(b)$

$$P(T,W)$$

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

$$P(W = s|T = c) = \frac{P(W = s, T = c)}{P(T = c)} = \frac{0.2}{0.5} = 0.4$$

$$= P(W = s, T = c) + P(W = r, T = c)$$
$$= 0.2 + 0.3 \ = 0.5$$

# Conditional Distributions

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

Conditional Distributions

$$P(W|T = hot)$$

| W | P |
|------|-----|
| sun | 0.8 |
| rain | 0.2 |

$P(W|T)$

$$P(W|T = cold)$$

| W | P |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

Joint Distribution

$$P(T, W)$$

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

$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)}$$

# Homework: Conditional Distribution

$P(X, Y)$

| X | Y | P |
|---|---|---|
| +x | +y | 0.2 |
| +x | -y | 0.3 |
| -x | +y | 0.4 |
| -x | -y | 0.1 |

- P(+x | +y) ?

- P(-x | +y) ?

- P(-y | +x) ?

# 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)

$P(T, W)$

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

**Select** $\longrightarrow$

$P(T, r)$

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

**Normalize** $\longrightarrow$

$P(T|r)$

| T | P |
|------|------|
| hot | 0.25 |
| cold | 0.75 |

  - Why does this work? Sum of selection is P(evidence)!  (P(r), here)

$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x_1, x_2)}{\sum_{x_1} P(x_1, x_2)}$$

# Normalization Trick

$$P(T, W)$$

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

$$P(W = s | T = c) = \frac{P(W = s, T = c)}{P(T = c)}$$
$$= \frac{P(W = s, T = c)}{P(W = s, T = c) + P(W = r, T = c)}$$
$$= \frac{0.2}{0.2 + 0.3} = 0.4$$

$$P(W = r | T = c) = \frac{P(W = r, T = c)}{P(T = c)}$$
$$= \frac{P(W = r, T = c)}{P(W = s, T = c) + P(W = r, T = c)}$$
$$= \frac{0.3}{0.2 + 0.3} = 0.6$$

$$P(W | T = c)$$

| W | P |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

# Normalization Trick

$$P(W = s | T = c) = \frac{P(W = s, T = c)}{P(T = c)}$$
$$= \frac{P(W = s, T = c)}{P(W = s, T = c) + P(W = r, T = c)}$$
$$= \frac{0.2}{0.2 + 0.3} = 0.4$$

$P(T, W)$

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

**SELECT** the joint probabilities matching the evidence

➡

$P(c, W)$

| T | W | P |
|------|------|-----|
| cold | sun | 0.2 |
| cold | rain | 0.3 |

**NORMALIZE** the selection (make it sum to one)

➡

$P(W | T = c)$

| W | P |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

$$P(W = r | T = c) = \frac{P(W = r, T = c)}{P(T = c)}$$
$$= \frac{P(W = r, T = c)}{P(W = s, T = c) + P(W = r, T = c)}$$
$$= \frac{0.3}{0.2 + 0.3} = 0.6$$

# Normalization Trick

$P(T, W)$

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

**SELECT** the joint probabilities matching the evidence

➡

$P(c, W)$

| T | W | P |
|------|------|-----|
| cold | sun | 0.2 |
| cold | rain | 0.3 |

**NORMALIZE** the selection (make it sum to one)

➡

$P(W|T = c)$

| W | P |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

- Why does this work? Sum of selection is P(evidence)! (P(T=c), here)

$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x_1, x_2)}{\sum_{x_1} P(x_1, x_2)}$$

# To Normalize

- (Dictionary) To bring or restore to a normal condition

  All entries sum to ONE

- Procedure:
  - Step 1: Compute Z = sum over all entries
  - Step 2: Divide every entry by Z

- Example 1

| W | P |
|------|-----|
| sun | 0.2 |
| rain | 0.3 |

Normalize →

Z = 0.5

| W | P |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

- Example 2

| T | W | P |
|------|------|----|
| hot | sun | 20 |
| hot | rain | 5 |
| cold | sun | 10 |
| cold | rain | 15 |

Normalize →

Z = 50

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

# Terminology

$c_i$

$y_j$       $n_{ij}$     $\Big\}\ r_j$

$x_i$

## Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

## Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

## Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

X value is given

# Probabilistic Inference

- Diagnosis
- Speech recognition
- Tracking objects
- Robot mapping
- Genetics
- Error correcting codes
- ... lots more!

# Probabilistic Inference

- **Probabilistic inference**: compute a desired probability from other known probabilities (e.g. conditional from joint)

- We generally compute conditional probabilities
  - P(on time | no reported accidents) = 0.90
  - These represent the agent's **beliefs** given the evidence

- Probabilities change with new evidence:
  - P(on time | no accidents, 5 a.m.) = 0.95
  - P(on time | no accidents, 5 a.m., raining) = 0.80
  - Observing new evidence causes beliefs to be updated

# Uncertainty

- **General situation:**

  - **Observed variables (evidence)**: Agent knows certain things about the state of the world (e.g., sensor readings or symptoms)

  - **Unobserved variables**: Agent needs to reason about other aspects (e.g. where an object is or what disease is present)

  - **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

# Inference by Enumeration

- P(sun)?

| S | T | W | P |
|---|---|---|---|
| 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

- 

- P(sun | winter)?

| S | T | W | P |
|---|---|---|---|
| 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

- 

- 

- <span style="color:red">P(sun | winter, hot)?</span>

| S | T | W | P |
|---|---|---|---|
| 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

- **General case:**
  - Evidence variables: $E_1 \ldots E_k = e_1 \ldots e_k$
  - Query* variable: $Q$
  - Hidden variables: $H_1 \ldots H_r$

  $X_1, X_2, \ldots X_n$

  All variables

- We want: $P(Q|e_1 \ldots e_k)$
- First, select the entries consistent with the evidence
- Second, sum out H to get joint of Query and evidence:

$$P(Q, e_1 \ldots e_k) = \sum_{h_1 \ldots h_r} \underbrace{P(Q, h_1 \ldots h_r, e_1 \ldots e_k)}_{X_1, X_2, \ldots X_n}$$

- Finally, normalize the remaining entries to conditionalize

# Problems with Enumeration

- Obvious problems:
  - Worst-case time complexity $O(d^n)$
  - Space complexity $O(d^n)$ to store the joint distribution

- Solutions
  - Better techniques
  - Better representation
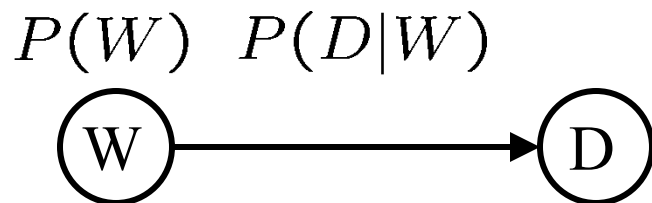  - Simplifying assumptions

# The Product Rule

- Sometimes have conditional distributions but want the joint

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

- Example:

$$P(D,W)$$

$$P(W) \quad P(D|W)$$

W → D

# The Chain Rule

- More generally, can always write any joint distribution as an incremental product of conditional distributions?

$$P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)$$

$$P(x_1, x_2, \ldots x_n) = \prod_i P(x_i|x_1 \ldots x_{i-1})$$

- Why is this always true?

# Bayes' Rule

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

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

That's my rule!

- Dividing, we get:

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

- Why is this at all helpful?
  - Lets us build a conditional from its reverse
  - Often one conditional is tricky but the other one is simple
  - Foundation of many systems we'll see later

- In the running for most important AI equation!

# Inference with Bayes' Rule

- **Example: Diagnostic probability from causal probability:**

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

- **Example:**

  - m is meningitis, s is stiff neck

$$P(s|m) = 0.8$$
$$P(m) = 0.0001$$
$$P(s) = 0.1$$

Example givens

$$P(m|s) = \frac{P(s|m)P(m)}{P(s)} = \frac{0.8 \times 0.0001}{0.1} = 0.0008$$

  - Note: posterior probability of meningitis still very small
  - Note: you should still get stiff necks checked out!  Why?

# Quiz: Bayes Rule

■ Given:

$P(W)$

| R | P |
|---|---|
| sun | 0.8 |
| rain | 0.2 |

$P(D|W)$

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

■ What is P(W | dry) ?

# Ghostbusters, Revisited

- Let's say we have two distributions:
  - Prior distribution over ghost location: P(G)
    - Let's say this is uniform
  - Sensor reading model: P(R | G)
    - Given: we know what our sensors do
    - R = reading color measured at (1,1)
    - E.g. P(R = yellow | G=(1,1)) = 0.1

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

- We can calculate the posterior distribution P(G|r) over ghost locations given a reading using Bayes' rule:

$$P(g|r) \propto P(r|g)P(g)$$

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

# Independence

- Two variables are *independent* if:

$$\forall x, y : P(x, y) = P(x)P(y)$$

  - This says that their joint distribution *factors* into a product two simpler distributions
  - Another form:

$$\forall x, y : P(x|y) = \phantom{xxxx}$$

  - We write: $X \perp\!\!\!\perp Y$

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

# Example: Independence?

$P_1(T, W)$

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

$P(T)$

| T | P |
|---|---|
| hot | 0.5 |
| cold | 0.5 |

$P(W)$

| W | P |
|---|---|
| sun | 0.6 |
| rain | 0.4 |

# Example: Independence

- N fair, independent coin flips:

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

$2^n \Big\{$

# Conditional Independence

- P(Toothache, Cavity, Catch)

- If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:

  - P(+catch | +toothache, +cavity) = P(+catch | +cavity)

- The same independence holds if I don't have a cavity:

  - P(+catch | +toothache, -cavity) = P(+catch| −cavity)

- Catch is *conditionally independent* of Toothache given Cavity:

  - P(Catch | Toothache, Cavity) = P(Catch | Cavity)

- Equivalent statements:

  - P(Toothache | Catch , Cavity) = P(Toothache | Cavity)

  - P(Toothache, Catch | Cavity) = P(Toothache | Cavity) P(Catch | Cavity)

  - One can be derived from the other easily

# Conditional Independence

- Unconditional (absolute) independence very rare (why?)

- *Conditional independence* is our most basic and robust form of knowledge about uncertain environments:

$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$$
$$\forall x, y, z : P(x|z, y) = P(x|z)$$

$$X \perp\!\!\!\perp Y | Z$$

- What about this domain:
  - Traffic
  - Umbrella
  - Raining

# Probability Summary

- **Conditional probability** $\qquad P(x|y) = \dfrac{P(x,y)}{P(y)}$

- **Product rule** $\qquad\qquad P(x,y) = P(x|y)P(y)$

- **Chain rule**
$$P(X_1, X_2, \ldots X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)\ldots$$
$$= \prod_{i=1}^{n} P(X_i|X_1, \ldots, X_{i-1})$$

- **X, Y independent if and only if:** $\quad \forall x, y : P(x,y) = P(x)P(y)$

- **X and Y are conditionally independent given Z if and only if:** $\quad X \perp\!\!\!\perp Y | Z$
$$\forall x, y, z : P(x,y|z) = P(x|z)P(y|z)$$