



Probabilistic Reasoning So far, mostly, we've done deterministic problems. Image: www.instructables.com/id/How-to-Win-a-Chess-Game-in-2-Moved

Probabilistic Reasoning

- We don't (can't!) know everything about most problems.
- Most problems are not:
 - Deterministic
 - Fully observable
- Or, we can't calculate everything.
 - Continuous problem spaces
- Probability lets us understand, quantify, and work with this uncertainty.



Probability

- World: The complete set of possible states
- Random variables: Problem aspects that take a value
 - "The number of blue squares we have pulled," B
 - "The combined value of two dice we rolled," ${\cal C}$
- Event: Something that happens
- Sample Space: All the things (outcomes) that could happen in some set of circumstances
 - Pull 2 squares from envelope A: what is the sample space?
 - How about envelope B?
- World, redux: A complete assignment of values to variables



<section-header><section-header><section-header><list-item><list-item>

Boolean Random Variables

- We will start with the simplest type of random variables: Booleans
- Take the values true or false
- Think of the event as occurring or not occurring
- Examples (Let A be a Boolean random variable):
 - A = Getting a head on a coin flip
 - A = It will rain today

8

Basic Probability

- Each P is a non-negative value in [0,1]
 - $P(\{1,1\}) = 1/36$
- Total probability of the sample space is 1
 - $P(\{1,1\}) + P(\{1,2\}) + P(\{1,3\}) + \dots + P(\{6,6\}) = 1$
- For mutually exclusive events, the probability for at least one of them is the sum of their individual probabilities
 - $P(sunny) \lor P(cloudy) = P(sunny) + P(cloudy)$
- Where do these numbers come from?
 - Experimental probability: Based on frequency of past events
 - Subjective probability: Based on expert assessment



The Joint Probability Distribution

- Joint probabilities can be between any number of variables
 - E.g. P(A = true, B = true, C = true)
- For each combination of variables, we need to say how probable that combination is
- The probabilities of these combinations need to sum to 1

| Α | B | С | P(A,B,C) |
|-------|-------|-------|-----------|
| false | false | false | 0.1 |
| false | false | true | 0.2 |
| false | true | false | 0.05 |
| false | true | true | 0.05 |
| true | false | false | 0.3 |
| true | false | true | 0.1 |
| true | true | false | 0.05 |
| true | true | true | 0.15 |
| | | | |
| | | | Sums to 1 |

10

The Joint Probability Distribution

- Once you have the joint probability distribution, you can calculate any probability involving A, B, and C
- Note: May need to use marginalization and Bayes rule
- Examples:
 - P(A=true) = sum of P(A,B,C) in rows with A=true
 - P(A=true, B = true | C=true) =
 P(A = true, B = true, C = true) / P(C = true)

| Α | B | С | P(A,B,C) |
|-------|-------|-------|-----------------|
| false | false | false | 0.1 |
| false | false | true | 0.2 |
| false | true | false | 0.05 |
| false | true | true | 0.05 |
| true | false | false | 0.3 |
| true | false | true | 0.1 |
| true | true | false | 0.05 |
| true | true | true | 0.15 |











Probability Distributions

- A distribution is the probabilities of **all possible values** of a random variable
- Ex: weather can be sunny, rainy, cloudy, or snowy
 - P(Weather = sun) = 0.6
 - P(Weather = rain) = 0.1
 - P(Weather = cloud) = 0.29
 - P(Weather = snow) = 0.01
 - $P(Weather) = \langle 0.6, 0.1, 0.29, 0.01 \rangle$ \leftarrow shortcut
- P(Weather): probability distribution on Weather



Probability Theory: Definitions

- **Conditional distribution** gives the probabilities of a variable, dependent on the values of the other variables
 - P(XIY) = "Probability of X happening, given that Y happens (or didn't happen)"
- Probability of some effect, given that we know cause(s)
 - (Technically, we only know b is correlated, not causal)
 - Example: P(*alarm* | *burglary*) = "Probability of *alarm*, given *burglary*"
- Computing it:

•
$$P(a \mid b) = \frac{P(a \land b)}{P(b)}$$

• P(b): **normalizing constant** (later we'll call this alpha α or rho ϱ)

Probability Theory: Definitions

• Product rule:

•

• $P(a \land b) = P(a \mid b) P(b)$

| Marginalizing | (summing out): |
|---------------|----------------|

- Finding distribution over one or a subset of variables
- Marginal probability of B summed over all alarm states:
- $P(B) = \Sigma_a P(B, a)$
 - P(B) = sum of P(B, a) for all possible values of A
- Conditioning over a subset of variables:
 - $P(B) = \Sigma_a P(B \mid a) P(a)$



| l et's Try It | | alarm | ¬ alarm |
|--|------------|-----------------|---------|
| | burgla | ry 0.09 | 0.01 |
| Cond'l probability | ¬ burgla | n ry 0.1 | 0.8 |
| • Cond i probability | | | |
| P(effect, cause[s]) | | | |
| P(a b) = P(a ∧ b) / P(b) | | | |
| P(b): normalizing constant (1/α) | | | |
| Product rule: | | | |
| P(a ∧ b) = P(a b) P(b) | | | |
| Marginalizing: | | | |
| P(B) = Σ_aP(B, a) | | | |
| • $P(B) = \Sigma_a P(B \mid a) P(a)$ (conditioning) | • $P(A) =$ | ? | |
| | | | |
| | | | |

| | alarm | ¬ alarm |
|-------------------|-------|---------|
| burglary | 0.09 | 0.01 |
| ¬ burglary | 0.1 | 0.8 |

Marginalizing

- Marginalization: how to safely ignore variables.
- Two-variable example (A and B).
- If we know P(A=a,B=b) for *all* values of *a* and *b*:
- $P(B=b)=\sum_{a}P(A=a,B=b).$
- Here we "marginalized out" the variable A.
- Takes variable(s) in a out of consideration

20

Marginalizing

- Marginal probability: the probability of an event occurring, regardless of other events (unlike conditional probability)
 - To get there we have to marginalize
- Marginalizing (summing out):
 - Finding distribution over one or a subset of variables
 - Marginal probability of B summed over all alarm states:
 - $P(B) = \Sigma_a P(B, a)$
- Takes variable(s) in *a* out of consideration





Exercise: Inference from the joint

- Queries:
 - What is the prior probability of *smart*?
 - What is the prior probability of *study*?
 - What is the conditional probability of *prepared*, given *study* and *smart*?

| P(smart ∧ | smart | | -smart | |
|---------------|-------|--------|--------|--------|
| study ∧ prep) | study | ¬study | study | ¬study |
| prepared | .432 | .16 | .084 | .008 |
| -prepared | .048 | .16 | .036 | .072 |

P(*smart*) = .432 + .16 + .048 + .16 = **0.8**

24

Exercise: Inference from the joint

- Queries:
 - What is the prior probability of *smart*?
 - What is the prior probability of *study*?
 - What is the conditional probability of *prepared*, given *study* and *smart*?

| P(smart ∧ | SM | art | <i>¬smart</i> | |
|---------------|-------|--------|---------------|--------|
| study ∧ prep) | study | ¬study | study | ¬study |
| prepared | .432 | .16 | .084 | .008 |
| ¬prepared | .048 | .16 | .036 | .072 |

| E | Exercise: Inference from the joint | | | | | | | |
|---|--|--|---------------|-------|--------|--------------------|--------|---|
| | Queries: What is the prior probability of <i>smart</i>? What is the prior probability of <i>study</i>? What is the conditional probability of <i>prepared</i>, given <i>study</i> and <i>smart</i>? | | | | | | | |
| | | | P(smart A | sn | art | ¬ <i>SN</i> | nart |] |
| | | | study ∧ prep) | study | ¬study | study | ¬study | |
| | | | prepared | .432 | .16 | .084 | .008 | |
| | | | -prepared | .048 | .16 | .036 | .072 | |
| | P(prep smart,study) = P(prep, smart, study)/P(smart, study) = .432 / (.432 + .048) = 0.9 | | | | | | | |

The Problem with the Joint Distribution

- Lots of entries in the table!
- For k Boolean random variables, you need a table of size 2k
- How do we use fewer numbers?
- Need independence

| Α | B | С | P(A,B,C) |
|-------|-------|-------|----------|
| false | false | false | 0.1 |
| false | false | true | 0.2 |
| false | true | false | 0.05 |
| false | true | true | 0.05 |
| true | false | false | 0.3 |
| true | false | true | 0.1 |
| true | true | false | 0.05 |
| true | true | true | 0.15 |

Independence

Variables A and B are independent if any of the following hold:

- P(A,B) = P(A) P(B)
- P(A | B) = P(A)
- P(B | A) = P(B)

This says that knowing the outcome of A does not tell me anything new about the outcome of B.

29

Independence

How is independence useful?

- Suppose you have n coin flips and you want to calculate the joint distribution $P(C_1, ..., C_n)$
- If the coin flips are not independent, you need 2ⁿ values in the table
- If the coin flips are independent, then

$$P(C_1,...,C_n) = \prod_{i=1}^n P(C_i)$$

Each $P(C_i)$ table has 2 entries and there are *n* of them for a total of 2*n* values



Independence Example

- {moon-phase, light-level} II {burglary, alarm, earthquake}
 - Unless maybe burglaries increase in low light?
 - {*light-level*} # {*burglary*}
 - But, if we know the light level, *moon-phase* \bot *burglary*
 - Once we're burglarized, light level doesn't affect whether the alarm goes off; {light-level} 1 {alarm}

• We need:

- 1. A more complex notion of independence
- 2. Methods for reasoning about these kinds of (common) relationships

Exercise: Independence

- Is *smart* independent of *study*?
 - Is P(*smart* | *study*) = P(*smart*) ?

• Is *prepared* independent of *study*?

• Is $P(prep \mid study) = P(prep)$?

| P(smart ∧ | smart | | <i>¬smart</i> | |
|---------------|-------|--------|---------------|--------|
| study ∧ prep) | study | ¬study | study | ¬study |
| prepared | .432 | .16 | .084 | .008 |
| ¬prepared | .048 | .16 | .036 | .072 |

33

Exercise: Independence

- Is *smart* independent of *study*?
 - $P(smart \mid study) = P(smart)$
- Is *prepared* independent of *study*?
 - $P(prep \mid study) = P(prep)$

| Smart | Study | | |
|-------|-------|---------------|-------|
| t | t | 0.432 + 0.48 | 0.480 |
| t | f | 0.16 + 0.16 | 0.32 |
| f | t | 0.084 + 0.008 | 0.092 |
| f | f | 0.036 + 0.72 | 0.756 |

| P(smart ∧ | smart | | <i>¬smart</i> | |
|---------------|-------|--------|---------------|--------|
| study ∧ prep) | study | ¬study | study | ¬study |
| prepared | .432 | .16 | .084 | .008 |
| ¬prepared | .048 | .16 | .036 | .072 |

Exercise: Independence

- Is $P(smart \mid study) = P(smart)$?
- Is P(*smart* | *study*) = P(*smart*, *study*) / P(*study*) ?
- 0.8 = (.432 + .048) / .6
- 0.8 = 0.8

| | | | Smart | Study | | | | |
|-------------------------|-------|--------|---------------|--------|----------|------------|---------------|-------|
| P(smart∧ study∧prep) | smart | | <i>¬smart</i> | | 5111a1 t | 5tuuy + | 0 422 ± 0 48 | 0.480 |
| | study | ¬study | study | ¬study | t | f | 0.432 + 0.48 | 0.480 |
| prepared | .432 | .16 | .084 | .008 | f | t | 0.084 + 0.008 | 0.092 |
| -prepared | .048 | .16 | .036 | .072 | f | f | 0.036 + 0.72 | 0.756 |
| | | | | | • | | | |

35

Conditional Probabilities

- Describes dependent events
 - Affect each other in some way
- Typical in the real world
- If we know some event has occurred, what does that tell us about the likelihood of another event?

Conditional Independence

- moon-phase and burglary are conditionally independent given light-level
 - That is, $M \perp B$ if we already know L
- Conditional independence is:
 - Weaker than absolute independence
 - Useful in decomposing full joint probability distributions



Conditional Independence

- Absolute independence: $A \perp B$, if:
 - $P(A \land B) = P(A) P(B)$
 - Equivalently, P(A) = P(A | B) and P(B) = P(B | A)
- A and B are conditionally independent given C if:
 - $P(A \land B \mid C) = P(A \mid C) P(B \mid C)$
- This lets us decompose the joint distribution:
 - $P(A \land B \land C) = P(A \mid C) P(B \mid C) P(C)$
- What does this mean?

Exercise: Conditional Independence

- Queries:
 - Is *smart* conditionally independent of *prepared*, given *study*?
 - Is *study* conditionally independent of *prepared*, given *smart*?

| P(smart ∧ | SM | art | <i>¬smart</i> | | |
|---------------|-------|--------|---------------|--------|--|
| study ∧ prep) | study | ¬study | study | ¬study | |
| prepared | .432 | .16 | .084 | .008 | |
| ¬prepared | .048 | .16 | .036 | .072 | |

39

Probability

- Worlds, random variables, events, sample space
- Joint probabilities of multiple connected variables
- Conditional probabilities of a variable, given another variable(s)
- Marginalizing out unwanted variables
- Inference from the joint probability

The big idea: figuring out the probability of variable(s) taking certain value(s)

Bayes' Rule

- Derive the probability of some event, given another event
 - Assumption of attribute independence (AKA the Naïve assumption)
 - Naïve Bayes assumes that all *attributes* are independent.
- Also the basis of modern machine learning
- Bayes' rule is derived from the product rule

$$P(Y \mid X) = \frac{P(X \mid Y) P(Y)}{P(X)}$$

R&N 495

41

Bayes' Rule

- $P(Y \mid X) = P(X \mid Y) P(Y) / P(X)$
- Often useful for diagnosis.
- If we have:
 - X = (observable) effects, e.g., symptoms
 - *Y* = (hidden) causes, e.g., illnesses
 - A model for how causes lead to effects: P(X | Y)
 - Prior beliefs about frequency of occurrence of effects: P(Y)
- We can reason from effects to causes: P(Y | X)

Naïve Bayes Algorithm

- Estimate the probability of each class:
 - Compute the posterior probability (Bayes rule)

$$P(c_i \mid D) = \frac{P(c_i)P(D \mid c_i)}{P(D)}$$

- Choose the class with the highest probability
- Assumption of attribute independency (Naïve assumption): Naïve Bayes assumes that all of the attributes are independent.



Simple Bayesian Diagnostic Reasoning

- We know:
 - Evidence / manifestations: $E_1, \ldots \, E_m$
 - Hypotheses / disorders: $H_1, \ldots H_n$
 - $E_{j} \mbox{ and } H_{i} \mbox{ are binary; hypotheses are mutually exclusive (non-overlapping) and exhaustive (cover all possible cases)$
 - Conditional probabilities: $P(E_j \mid H_i), \, i=1, \, \ldots \, n; \, j=1, \, \ldots \, m$
- Cases (evidence for a particular instance): $E_1, ..., E_m$
- Goal: Find the hypothesis H_i with the highest posterior
 - $Max_i P(Hi | E_1, \dots, E_m)$

45

Priors

- Four values total here:
 - P(H | E) = (P(E | H) * P(H)) / P(E)
- P(H | E) what we want to compute
- Three we already know, called the priors
 - $P(E \mid H)$
 - P(H)
 - P(E)

(In ML later, we will use the training set to estimate the priors)

Bayesian Diagnostic Reasoning II

- Bayes' rule says that
 - $P(H_i | E_1, ..., E_m) = P(E_1, ..., E_m | H_i) P(H_i) / P(E_1, ..., E_m)$
- Assume each piece of evidence E_i is **conditionally independent** of the others, **given** a hypothesis H_i , then:
 - $P(E_1, ..., E_m | H_i) = \prod_{j=1}^{l} P(E_j | H_i)$
- If we only care about relative probabilities for the H_i , then we have:
 - $P(H_i | E_1, ..., E_m) = \alpha P(H_i) \prod_{j=1}^{l} P(E_j | H_i)$







A Bayesian Network

A Bayesian network is made up of:

I.A Directed Acyclic Graph



2.A set of tables for each node in the graph

| А | P(A) | Α | В | P(B A) | В | D | P(D B) | В | С | P(C B) |
|-------|------|-------|-------|--------|-------|-------|--------|-------|-------|--------|
| false | 0.6 | false | false | 0.01 | false | false | 0.02 | false | false | 0.4 |
| true | 0.4 | false | true | 0.99 | false | true | 0.98 | false | true | 0.6 |
| _ | | true | false | 0.7 | true | false | 0.05 | true | false | 0.9 |
| | | true | true | 0.3 | true | true | 0.95 | true | true | 0.1 |



A Bayesian Network

A Bayesian network is made up of:

I.A Directed Acyclic Graph



2.A set of tables for each node in the graph

| А | P(A) | Α | В | P(B A) | В | D | P(D B) | В | С | P(C B) |
|-------|------|-------|-------|--------|-------|-------|--------|-------|-------|--------|
| false | 0.6 | false | false | 0.01 | false | false | 0.02 | false | false | 0.4 |
| true | 0.4 | false | true | 0.99 | false | true | 0.98 | false | true | 0.6 |
| | | true | false | 0.7 | true | false | 0.05 | true | false | 0.9 |
| | | true | true | 0.3 | true | true | 0.95 | true | true | 0.1 |









<section-header><list-item><list-item><list-item><list-item><list-item>Example: Toothache • Random variables: • How's the weather? • Do you have a toothache? • Does the dentist's probe catch when she pokes your tooth? • Do you have a cavity? Weather Cavity Toothache Cavity Stiles derived from Mart E. Taylor, U. Alberte















Example BN

 We only specify P(A) etc., not P(¬A), since they have to sum to one



65

Bayesian Belief Networks (BNs) Making a Bayesian Network BN: BN = (DAG, CPD) DAG: directed acyclic graph (BN's structure) Nodes: random variables Typically binary or discrete Methods exist for continuous variables Arcs: indicate probabilistic dependencies between nodes Lack of link signifies conditional independence CPD: conditional probability distribution (BN's parameters) Conditional probabilities at each node, usually stored as a table (conditional probability table, or CPT)



Probabilities in BNs

- Bayes' nets implicitly encode joint distributions as a product of local conditional distributions.
- To see probability of a full assignment, multiply all the relevant conditionals together:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i \mid parents(X_i))$$

- Example: P(+cavity, +catch, ¬toothache) = ?
- This lets us reconstruct any entry of the full joint



























Summary

- Probability review
 - Distributions, conditional probability, marginalizing
 - Independence
 - Bayes' rule
- Bayes' nets (Bayesian Belief Networks)
 - Graphical notation
 - Conditional probability tables
 - Probability distributions
- Next time
 - Inference using Bayes' nets