

Bayes Nets

CPS 170

Ron Parr

Why Joint Distributions are Important

- Joint distributions gives $P(X_1 \dots X_n)$
- Classification/Diagnosis
 - Suppose X_1 =disease
 - $X_2 \dots X_n$ = symptoms
- Co-occurrence
 - Suppose X_3 =lung cancer
 - X_5 =smoking
- Rare event Detection
 - Suppose $X_1 \dots X_n$ = parameters of a credit card transaction
 - Call card holder if $P(X_1 \dots X_n)$ is below threshold?

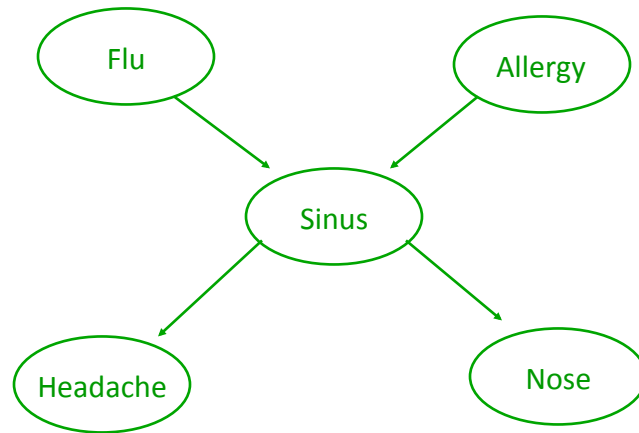
Modeling Joint Distributions

- To do this correctly, we need a full assignment of probabilities to all atomic events
- Unwieldy in general for discrete variables: n binary variables = 2^n atomic events
- Independence makes this tractable, but too strong (rarely holds)
- Conditional independence is a good compromise: Weaker than independence, but still has great potential to simplify things

Conditional Independence

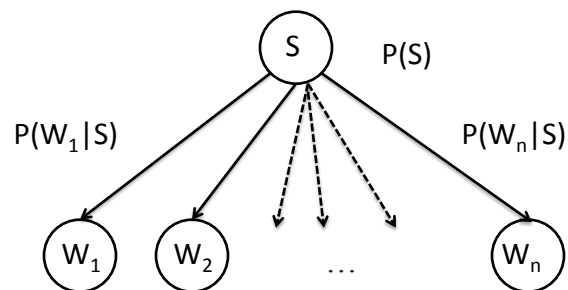
- Suppose we know the following:
 - The flu causes sinus inflammation
 - Allergies cause sinus inflammation
 - Sinus inflammation causes a runny nose
 - Sinus inflammation causes headaches
- How are these connected?

Graphical Structure



Knowing sinus separates the variables from each other.

Naïve Bayes Spam Filter



We will see later why this is a particularly convenient representation.
(Does it make a correct assumption?)

Conditional Independence

- We say that two variables, A and B, are conditionally independent given C if:
 - $P(A|BC) = P(A|C)$
 - $P(AB|C) = P(A|C)P(B|C)$
- How does this help?
- We **store only a conditional probability table** (CPT) of each variable given its parents
- Naïve Bayes (e.g. Spam Assassin) is a special case of this!

Notation Reminder

- $P(A|B)$ is a conditional prob. distribution
 - **It is a function!**
 - $P(A=\text{true}|B=\text{true})$, $P(A=\text{true}|B=\text{false})$, $P(A=\text{false}|B=\text{True})$, $P(A=\text{false}|B=\text{true})$
- $P(A|b)$ is a probability distribution, function
- $P(a|B)$ is a function, not a distribution
- $P(a|b)$ is a number

What is Bayes Net, More Formally

- A directed acyclic graph (DAG)
- Given parents, each variable is *independent of non-descendants*
- Joint probability decomposes:

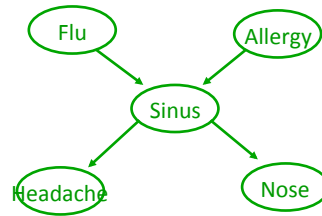
$$P(x_1 \dots x_n) = \prod_i P(x_i | \text{parents}(x_i))$$

- For each node X_i , store $P(X_i | \text{parents}(X_i))$
- Call this a Conditional Probability Table (CPT)
- CPT size is exponential in number of parents

Real Applications of Bayes Nets

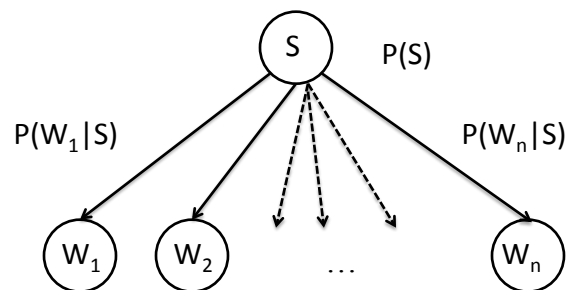
- Diagnosis of lymph node disease
- Used in Microsoft office and Windows
 - <http://research.microsoft.com/en-us/groups/mlas/>
- Used by robots to identify meteorites to study
- Study the human genome: Alex Hartemink et al.
- Many other applications...

Space Efficiency



- Entire joint distribution as 32 (31) entries
 - $P(H|S), P(N|S)$ have 4 (2)
 - $P(S|AF)$ has 8 (4)
 - $P(A)$ has 2 (1)
 - Total is 20 (10)
- This can require exponentially less space
- **Space problem is solved** for “most” problems

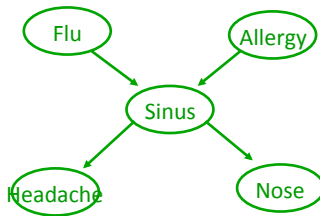
Naïve Bayes Space Efficiency



Entire Joint distribution has $2^{n+1} (2^{n+1}-1)$ numbers vs. $4n+2 (2n+1)$

Atomic Event Probabilities

$$P(x_1 \dots x_n) = \prod_i P(x_i \mid \text{parents}(x_i))$$



Note that this is guaranteed true if we construct net incrementally, so that for each new variable added, we connect all influencing variables as parents (prove it by induction)

(Non)Uniqueness of Bayes Nets

- You can always construct a valid Bayes net by inserting variables one at a time
- Order of adding variables can lead to different Bayesian networks for the same distribution
- Suppose A and B are independent, but C is a function of A and B
 - Add A, B, then C:
 - Add C, A, then B:

Doing Things the Hard Way

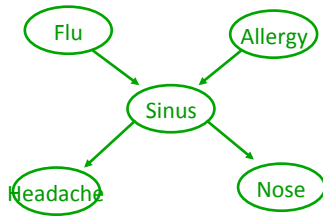
$$P(f|h) = \frac{P(fh)}{P(h)} = \frac{\sum_{SAN} P(fhSAN)}{\sum_{SANF} P(hSANF)} = \frac{\sum_{SAN} P(f)P(A)P(S|Af)P(h|S)P(N|S)}{\sum_{SANF} P(F)P(A)P(S|AF)P(h|S)P(N|S)}$$

defn. of conditional probability

marginalization

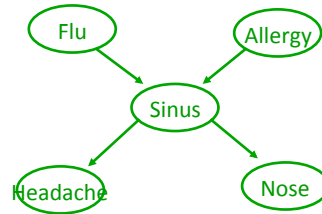
Doing this naïvely, we need to sum over all atomic events defined over these variables. There are exponentially many of these.

Working Smarter I



$$\begin{aligned}
 P(hSANF) &= \prod_x p(x | \text{parents}(x)) \\
 &= P(h|S)P(N|S)P(S|AF)P(A)P(F)
 \end{aligned}$$

Working Smarter II



$$\begin{aligned}
 P(h) &= \sum_{SANF} P(hSANF) \\
 &= \sum_{SANF} P(h|S)P(N|S)P(S|AF)P(A)P(F) \\
 &= \sum_{NS} P(h|S)P(N|S) \sum_{AF} P(S|AF)P(A)P(F) \\
 &= \sum_S P(h|S) \sum_N P(N|S) \sum_{AF} P(S|AF)P(A)P(F)
 \end{aligned}$$

Potential for exponential reduction in computation.

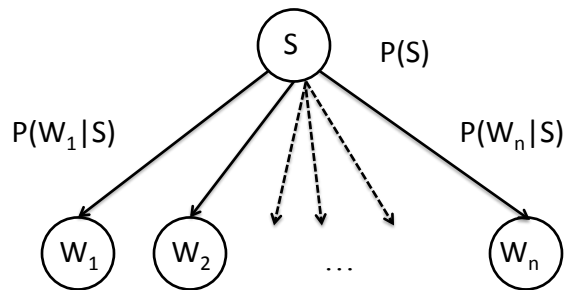
Computational Efficiency

$$\begin{aligned}
 \sum_{SANF} P(hSANF) &= \sum_{SANF} P(h|S)P(N|S)P(S|AF)P(A)P(F) \\
 &= \sum_S P(h|S) \sum_N P(N|S) \sum_{AF} P(S|AF)P(A)P(F)
 \end{aligned}$$

The distributive law allows us to decompose the sum.
AKA: Sum-product algorithm

Potential for an exponential reduction in computation costs.

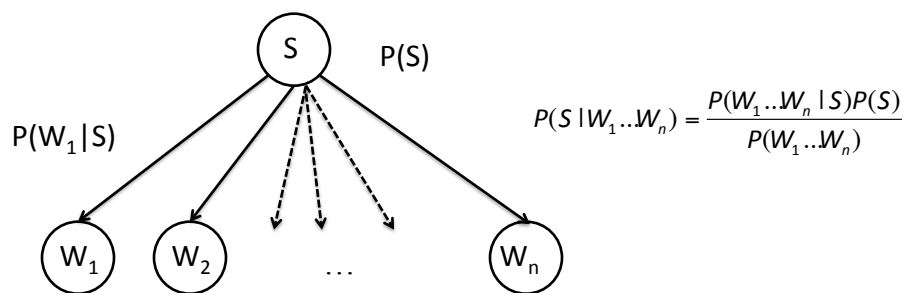
Naïve Bayes Efficiency



Given a set of words, we want to know which is larger: $P(s|W_1...W_n)$ or $P(\neg s|W_1...W_n)$.

Use Bayes Rule:
$$P(S|W_1...W_n) = \frac{P(W_1...W_n|S)P(S)}{P(W_1...W_n)}$$

Naïve Bayes Efficiency II



Observation 1: We can ignore $P(W_1...W_n)$

Observation 2: $P(S)$ is given

Observation 3: $P(W_1...W_n|S)$ is easy:
$$P(W_1...W_n|S) = \prod_{i=1}^n P(W_i|S)$$

Checkpoint

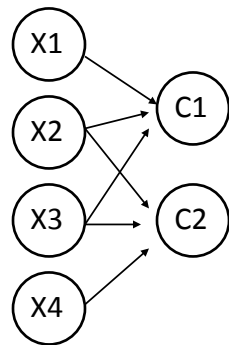
- BNs can give us an **exponential reduction** in the space required to represent a joint distribution.
- Storage is exponential in largest parent set.
- Claim: Parent sets are often reasonable.
- Claim: Inference cost is often reasonable.
- Question: Can we quantify relationship between structure and inference cost?

Now the Bad News...

- In full generality: Inference is NP-hard
- Decision problem: Is $P(X) > 0.5$?
- We reduce from 3SAT
- 3SAT variables map to BN variables
- Clauses become variables with the corresponding SAT variables as parents

Reduction

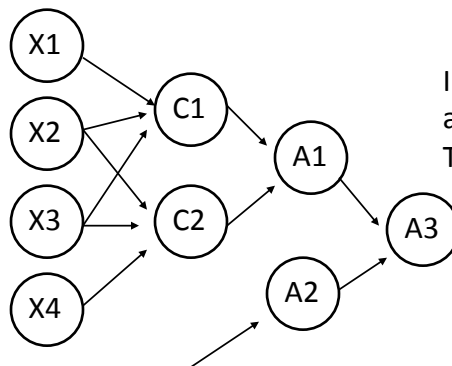
$$(\bar{X}_1 \vee X_2 \vee X_3) \wedge (\bar{X}_2 \vee X_3 \vee X_4) \wedge \dots$$



Problem: What if we have a large number of clauses?
How does this fit into our decision problem framework?

And Trees

We could make a single variable which is the AND of all of our clauses, but this would have CPT that is exponential in the number of clauses.



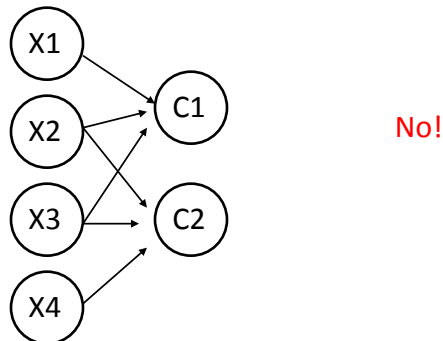
Implement as a tree of ANDs.
This is polynomial.

Checkpoint

- BNs can be very compact
- Worst case: Inference is intractable
- Hope that worst is case:
 - Avoidable (frequently, but no free lunch)
 - Easily characterized in some way

(Undirected) Trees

- Are the structures from our reduction trees?

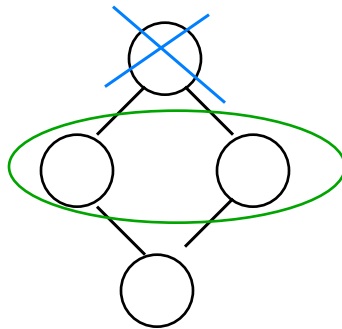


- They are a Directed Acyclic Graph (DAG)
- BNs are always DAGs (sometimes trees)

Clues in the Graphical Structure

- Q: How does graphical structure relate to our ability to push in summations over variables?
- A:
 - We relate summations to graph operations
 - Summing out a variable =
 - Removing node(s) from DAG
 - Creating new replacement node
 - Relate graph properties to computational efficiency

Variable Elimination

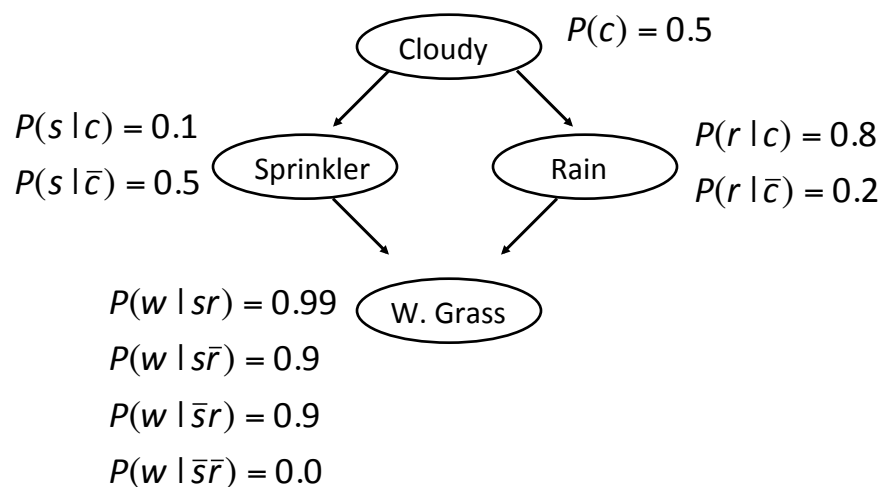


Recall that in variable elimination for CSPs, we eliminated variables and created new supervariables

BN Variable Elimination

- The same trick applies to Bayes nets
- Observation:
 - Every variable elimination ordering corresponds to a rearrangement of the summation in the marginalization computation
 - Variable elimination and sum-product are essentially the same algorithm

Another Example Network

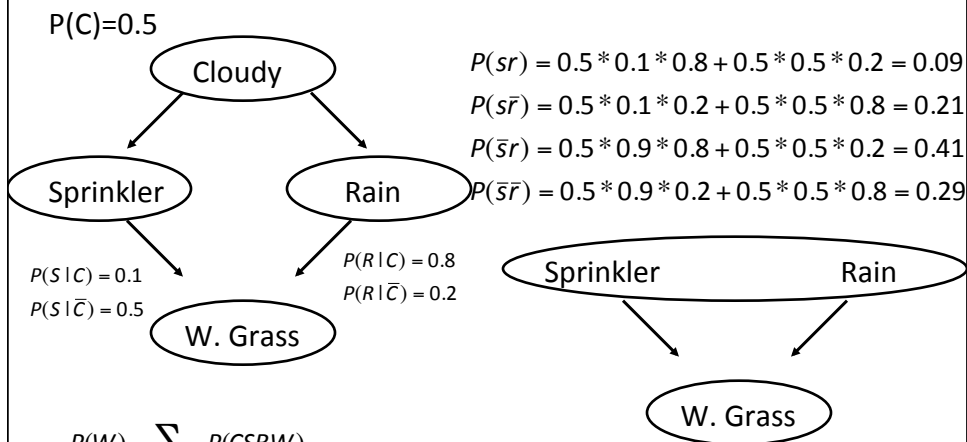


Marginal Probabilities

Suppose we want $P(W)$:

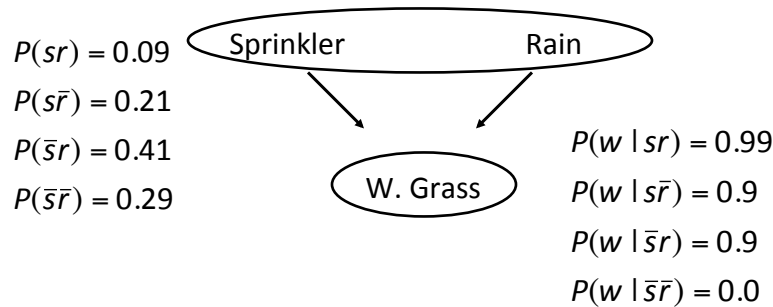
$$\begin{aligned}
 P(W) &= \sum_{CSR} P(CSRW) \\
 &= \sum_{CSR} P(C)P(S|C)P(R|C)P(W|RS) \\
 &= \sum_{SR} P(W|RS) \sum_C P(S|C)P(C)P(R|C)
 \end{aligned}$$

Eliminating Cloudy



$$\begin{aligned}
 P(W) &= \sum_{CSR} P(CSRW) \\
 &= \sum_{CSR} P(C)P(S|C)P(R|C)P(W|RS) \\
 &= \sum_{SR} P(W|RS) \sum_C P(S|C)P(C)P(R|C)
 \end{aligned}$$

Eliminating Sprinkler/Rain



$$\begin{aligned}
 P(w) &= \sum_{SR} P(w | RS)P(RS) \\
 &= 0.09 * 0.99 + 0.21 * 0.9 + 0.41 * 0.9 + 0.29 * 0 \\
 &= 0.6471
 \end{aligned}$$

Dealing With Evidence

Suppose we have observed that the grass is wet?
What is the probability that it has rained?

$$\begin{aligned}
 P(R | W) &= \alpha P(RW) \\
 &= \alpha \sum_{CS} P(CSRW) \\
 &= \alpha \sum_{CS} P(C)P(S | C)P(R | C)P(W | RS) \\
 &= \alpha \sum_C P(R | C)P(C) \sum_S P(S | C)P(W | RS)
 \end{aligned}$$

Is there a more clever way to deal with w?

Turning our Summation Trick into an Algorithm

- What happens when we “sum out” a variable?
 - All CPTs that reference this variable get pushed to the right of the summation
 - A new function defined over the union of these variables replaces these CPTs
- This is the Bayes net version of variable elimination from CSPs
- Analogous to Gaussian elimination in many ways but more expensive because operations are exponential in number of variables involved, rather than linear

The Variable Elimination Algorithm

```

Elim(bn, query)
If bn.vars = query
  return bn
Else
  x = pick_variable(bn)
  newbn.vars = bn.vars - x
  newbn.vars = newbn.vars - neighbors(x)
  newbn.vars = newbn.vars + newvar
  newbn.vars(newvar).function =

```

Can also sum
out variables
that are “hidden”

$$\longrightarrow \sum_X \prod_{Y \in X \cup \text{neighbors}(X)} \text{bn.vars}(Y).\text{function}$$

```

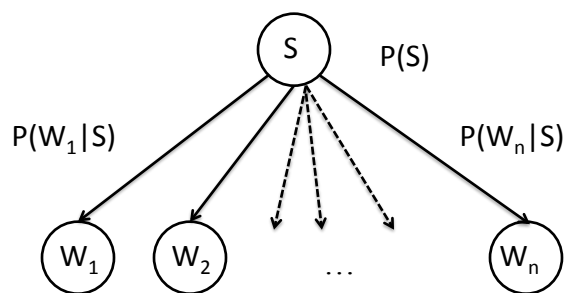
return(elim(newbn, query))

```

Efficiency of Variable Elimination

- Exponential in the largest domain size of new variables created (just as in CSPs)
- Equivalently: Exponential in largest function created by pushing in summations (sum-product algorithm)
- Linear for trees
- Almost linear for almost trees 😊

Naïve Bayes Efficiency



Another way to understand why Naïve Bayes is efficient:
It's a tree!

Facts About Variable Elimination

- Picking variables in optimal order is NP hard
- For some networks, there will be no elimination ordering that results in a poly time solution
(Must be the case unless $P=NP$)
- Polynomial for trees
- Need to get a little fancier if there are a large number of query variables or evidence variables

Beyond Variable Elimination

- Variable elimination must be rerun for every new query
- Possible to compile a Bayes net into a new data structure to make repeated queries more efficient
 - Recall that inference in trees is linear
 - Define a “cluster tree” where
 - Clusters = sets of original variables
 - Can infer original probs from cluster probs
- For networks w/o good elimination schemes
 - Sampling (discussed briefly)
 - Cutsets (discussed in text)
 - Variational methods (not covered in this class)
 - Loopy belief propagation (not covered in this class)

Sampling

- A Bayes net is an example of a **generative model** of a probability distribution
- Generative models allow one to generate samples from a distribution in a natural way
- Sampling algorithm:
 - While some variables are not sampled
 - Pick variable x with no unsampled parents
 - Assign this variable a value from $p(x | \text{parents}(x))$
 - Do this n times
 - Compute $P(a)$ by counting in what fraction a is true

Comments on Sampling

- Sampling is the easiest algorithm to implement
- Can compute marginal or conditional distributions by counting
- Not efficient in general
- Problem: How do we handle observed values?
 - Rejection sampling: Quit and start over when mismatches occur
 - Importance sampling: Use a reweighting trick to compensate for mismatches
- More clever approaches to sampling are possible

Summary of Algorithms for BNs

- Enumeration (consider all atomic events)
 - Exponential
 - Yuck!
- Sum-product/Variable elimination
 - Can be dramatically more efficient than enumeration
 - No guarantee of polynomial time
- Sampling
 - Easy to implement
 - May converge slowly in practice
- Active research area: Many other approaches used

Bayes Net Summary

- Bayes net = data structure for joint distribution
- Can give exponential reduction in storage
- Variable elimination:
 - simple, elegant method
 - efficient for many networks
- For some networks, must use approximation
- BNs are a major success story for modern AI
 - BNs do the “right” thing (no ugly approximations)
 - Exploit structure in problem to reduce storage/computation
 - Not always efficient, but inefficient cases are well understood
 - Work and used in practice