

Bayes Nets

CPS 271
Ron Parr

Modeling Distributions

- Suppose we knew $P(X_1 \dots X_n)$ for all features
 - Can answer any classification question optimally
 - Let $Y = X_i$
 - $P(Y | X_1 \dots X_n, X_i)$
 - Can answer many clustering type questions
 - $P(X_i X_j)$? (How often do two features co-occur)
 - $P(X_1 \dots X_n)$ (How typical is an instance?)
- To do correctly we need joint probability distribution
- Unwieldy for discrete variables
- Use *independence* to make this tractable

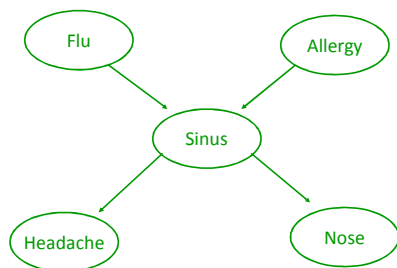
Where This Is Going

- Want: Some clever data structures and algorithms to circumvent the combinatorial explosion in the size of the joint distribution
- Note: BNs are *NOT* a learning method
- First: Understand how to use these data structures
- Relevance to machine learning:
 - Very useful to assume/have such structures
 - Learning of parameters
 - Learning of structure

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?

Causal Structure



Knowing sinus separates the variables from each other.

Conditional Independence

- We say that two variables, A and B, are conditionally independent given C if:
 - $P(A | BC) = P(A | C)$
- 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

Getting More Formal

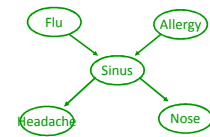
- What is a Bayes net?
 - A directed acyclic graph (DAG)
 - Given the parents, each variable is independent of non-descendants
 - Joint probability decomposes:

$$P(x_1 \dots x_n) = \prod_i P(x_i | \text{parents}(x_i))$$
 - For each node X_i , store $P(X_i | \text{parents}(X_i))$
 - Represent as table called a CPT

Real Applications of Bayes Nets

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

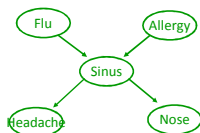
Space Efficiency



- Entire joint 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

Atomic Event Probabilities

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



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)

Doing Things the Hard Way

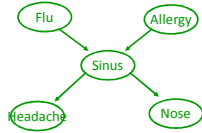
$$P(f|h) = \frac{P(fh)}{P(h)} = \frac{\sum_{SAN} P(fhSAN)}{\sum_{SANF} P(hSANF)}$$

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.

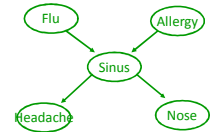
Working Smarter I



$$P(hSANF) = \prod_x p(x \mid \text{parents}(x))$$

$$= P(h \mid S)P(N \mid S)P(S \mid AF)P(A)P(F)$$

Working Smarter II



$$P(h) = \sum_{SANF} P(hSANF)$$

$$= \sum_{SANF} P(h \mid S)P(N \mid S)P(S \mid AF)P(A)P(F)$$

$$= \sum_{NS} P(h \mid S)P(N \mid S) \sum_{AF} P(S \mid AF)P(A)P(F)$$

$$= \sum_S P(h \mid S) \sum_N P(N \mid S) \sum_{AF} P(S \mid AF)P(A)P(F)$$

Potential for exponential reduction in computation.

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?

Computational Efficiency

$$\sum_{SANF} P(hSANF) = \sum_{SANF} P(h \mid S)P(N \mid S)P(S \mid AF)P(A)P(F)$$

$$= \sum_S P(h \mid S) \sum_N P(N \mid S) \sum_{AF} P(S \mid AF)P(A)P(F)$$

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

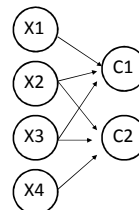
Potential for an exponential reduction in computation costs.

Now the Bad News...

- In full generality: Inference is NP-hard
- Decision problem: Is $P(X) > 0$?
- We reduce from 3SAT
- 3SAT variables map to BN variables
- Clauses become variables with the corresponding SAT variables as parents

Reduction

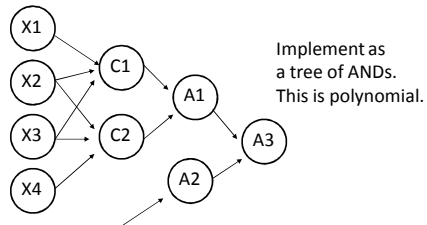
$$(\bar{X}_1 \vee X_2 \vee X_3) \wedge (\bar{X}_2 \vee X_3 \vee X_4) \wedge \dots$$



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.



Is BN Inference NP Complete?

- Can show that BN inference is #P hard
- #P is counting the number of satisfying assignments
- Idea: Assign variables uniform probability
- Probability of conjunction of clauses tells us how many assignments are satisfying

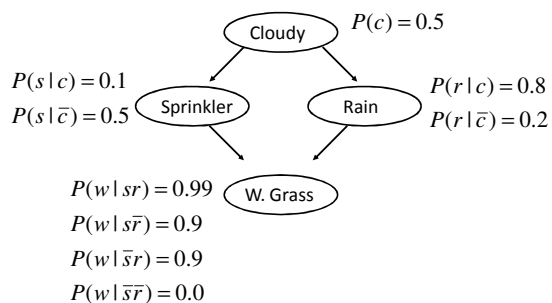
Checkpoint

- BNs can be very compact
- Worst case: Inference is intractable
- Hope that worst is case:
 - Avoidable
 - Easily characterized in some way

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

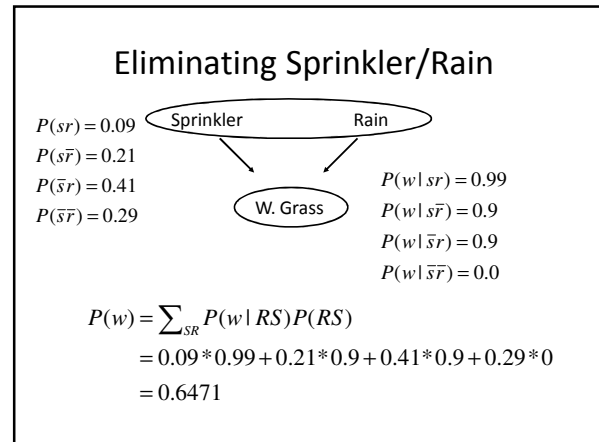
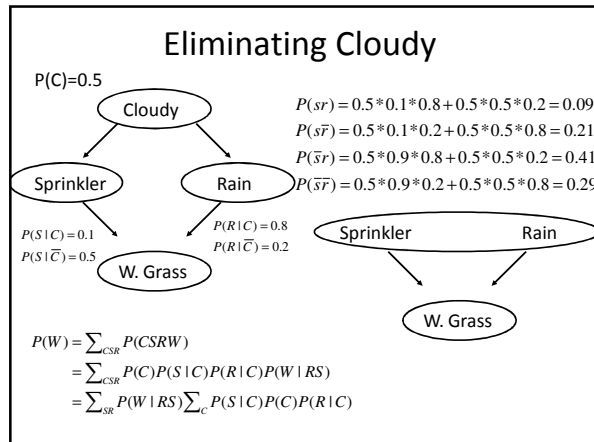
Another Example Network



Marginal Probabilities

Suppose we want $P(W)$:

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



Dealing With Evidence

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

$$P(R|W) = \alpha P(RW)$$

$$= \alpha \sum_{CS} P(CSRW)$$

$$= \alpha \sum_{CS} P(C)P(S|C)P(R|C)P(W|RS)$$

$$= \alpha \sum_C P(R|C)P(C) \sum_S P(S|C)P(W|RS)$$

Is there a more clever way to deal with w?

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

  return(elim(newbn, query))
  
```

Efficiency of Variable Elimination

- Exponential in the largest domain size of new variables created
- Equivalently: Exponential in largest function created by pushing in summations (sum-product algorithm)
- Linear for trees
- Almost linear for almost trees ☺
- (See examples on board...)

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
 - Note 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
 - Variational methods

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

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
- Q: Why is this interesting for machine learning?
 - A1: Very useful data structure!
 - A2: Often necessary to assume structure (even if it isn't quite right)
 - A3: Learning/discovering structure can be very useful