

We wish to learn both the parameters and structure of a Bayesian Network.

Parameters

Consider the disease model from the previous lecture:

2

Cases 1 and 3 have incomplete data, whereas 2 had complete.
If every example is complete, the data set is called complete.
If complete, the maximum likelihood parameters are unique.
We find the likelihood of a parameter set like so:

4

Param $\{S_1\} \rightarrow BN_1 \rightarrow Pr_1(.)$

Param $\{S_2\} \rightarrow BN_2 \rightarrow Pr_2(.)$

.

.

.

We pick the set of parameters that maximizes the probability,
so score $S_i = \prod Pr_i(e_i)$ is the likelihood of parameters

Parameter Estimation With Complete Data

7

What do I need to optimize to choose between these three structures?

We have learned many of these algorithms, so we won't discuss in detail, but:

1. local search methods (approximate methods)
→ we transform the structure looking for a better score
We must thus specify movement within the neighborhood structure: → legal operations:
add/remove/reverse edge
2. systemic search methods (exact methods)
→ A* is a good example

Why can't we just use maximum likelihood? We face **overfitting**.

Say we are using likelihood to compare; then $C > B > A$.

We will thus end up with a complete DAG, no matter what.

Thus we need a model that balances structure complexity with likelihood.

Why? Consider:

The data is clearly a linear fit, but if we estimate to the fourth degree, we get bad data!
There is no fixed answer to the problem; we clearly need some function of the form $f(\text{likelihood}) - g(\text{complexity})$.
A common one is called MDK, but we need not discuss details.

Model-Oriented Vs. Query-Oriented Learning

This can also be referred to as (unsupervised vs supervised) or (unlabeled vs labeled).
We have done model learning now; we move on to query based.

A model-based approach might give us the following:

)

Say we promise only to infer upward:

Then we don't 'model' the system and instead prepare for a specific 'query'.

Labeling of correct responses is done by humans, and this is done to 'train'.

If we want to talk about supervised learning, we need to introduce

Arithmetic Circuits

The (+) symbols represent OR gates
The (*) symbols represent AND gates
The θ are the parameters
The λ are the evidence

When given a query, we can perform weighted model counting on the equivalent arithmetic circuit.

If A=True, $\lambda_A = 1$; if A=False, $\lambda_{\neg A} = 1$; if unknown, both are 1

Thus we can evaluate this in linear time; unfortunately converting to an arithmetic circuit is $O(nd^w)$

BUT what do we do if we don't have the parameters? We use the labeled data:

$\text{![]}(https : //paper - attachments.dropbox.com/s_E29353D8DE6A7F32069419A77A11E9F7A8BD49E718AC081FD5F77701FB86FF68_15909113848drawing + 7.jpg)$

Cross Entropy gives us a measure of the disagreement between the two.

Therefore, we often seek to minimize it across a structure.

We generally use it to perform gradient descent (nth dimensional hill climbing).

Consider the example of recognizing shapes:

)

Though pixel is functional in general, we allow it to be probabilistic to permit noise.

It will thus be inferred to be very close to 0/1.

A lot of heights will be zeroed as well depending on the column.

We can thus see that we have a lot of **background knowledge**
for which we could simply substitute and reduce cases.

The simplest case of this is $A \iff B$,