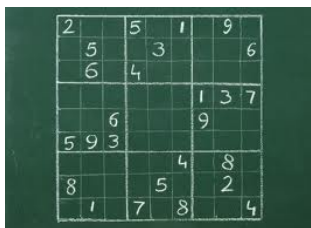


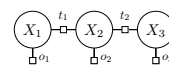


Lecture 6.2: CSPs III



Review: Object tracking

Factor graph (chain-structured):



- Variables X_i : location of object at time i
- Observation factors $o_i(x_i)$: noisy information compatible with position
- Transition factors $t_i(x_i, x_{i+1})$: object positions can't change too much



Review: Inference in CSPs

Algorithms for max-weight assignments in factor graphs:

(1) Extend partial assignments:

- Backtracking search: exact, exponential time
- Beam search: approximate, linear time

(2) Modify complete assignments:

- Iterated conditional modes: approximate, deterministic

- Last time, we started looking into local search procedures. In contrast to backtracking search or beam search, which build partial assignments one variable at a time, local search starts with an arbitrary complete assignment and iteratively improves it.

Review: Iterated conditional modes (ICM)

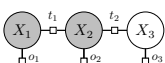
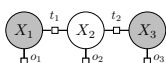
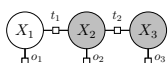


Algorithm: iterated conditional modes (ICM)

Initialize x to a random complete assignment

Loop through $i = 1, \dots, n$ until convergence:

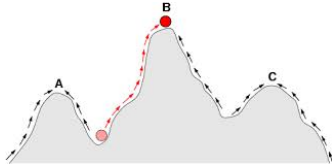
 Compute weight of $x_v = x \cup \{X_i : v\}$ for each v
 $x \leftarrow x_v$ with highest weight



- The first local search algorithm we looked at was Iterated conditional modes (ICM). ICM starts with a random complete assignment, then loops through each variable X_i . At each step, ICM updates the value of X_i to the best value $v \in \text{Domain}_i$, assuming all other variables remain fixed.

Review: Iterated conditional modes (ICM)

- Can get stuck in **local optima**
- Not guaranteed to find optimal assignment!



- Recall that ICM can get stuck in local optima, where there is a better assignment elsewhere, but all the one variable changes result in a lower weight assignment.



Roadmap

Gibbs Sampling

Conditioning

Elimination



Inference in CSPs

Algorithms for max-weight assignments in factor graphs:

(1) Extend partial assignments:

- Backtracking search: exact, exponential time
- Beam search: approximate, linear time

(2) Modify complete assignments:

- Iterated conditional modes: approximate, deterministic
- **Gibbs sampling**: approximate, randomized

Gibbs sampling

Sometimes, need to go downhill to go uphill...



Key idea: randomness

Sample an assignment with probability proportional to its weight.

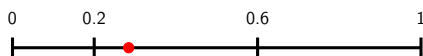


Example: Gibbs sampling

$\text{Weight}(x \cup \{X_2 : 0\}) = 1$ prob. 0.2

$\text{Weight}(x \cup \{X_2 : 1\}) = 2$ prob. 0.4

$\text{Weight}(x \cup \{X_2 : 2\}) = 2$ prob. 0.4



- In reinforcement learning, we also had a problem where if we explore by using a greedy policy (always choosing the best action according to our current estimate of the Q function), then we were doomed to get stuck. There, we used **randomness** via epsilon-greedy to get out of local optima.
- Here, we will do the same, but using a slightly more sophisticated form of randomness. The idea is **Gibbs sampling**, a method originally designed for using Markov chains to sample from a distribution over assignments. We will return to that original use later, but for now, we are going to repurpose it for the problem of finding the maximum weight assignment.

Gibbs sampling

[demo: gibbsSampling()]



Algorithm: Gibbs sampling

Initialize x to a random complete assignment
Loop through $i = 1, \dots, n$ until convergence:
 Compute weight of $x_v = x \cup \{X_i : v\}$ for each v
 Choose $x \leftarrow x_v$ with probability prop. to its weight

Can escape from local optima (not always easy though)!

- The form of the algorithm is identical to ICM. The only difference is that rather than taking the assignment $x \cup \{X_i : v\}$ with the maximum weight, we choose the assignment with probability proportional to its weight.
- In this way, even if an assignment has lower weight, we won't completely rule it out, but just choose it with lower probability. Of course if an assignment has zero weight, we will choose it with probability zero (which is to say, never).
- Randomness is not a panacea and often Gibbs sampling can get ensnared in local optima just as much as ICM. In theory, under the assumption that we could move from the initial assignment and the maximum weight assignment with non-zero probability, Gibbs sampling will move there eventually (but it could take exponential time in the worst case).
- Advanced: just as beam search is greedy search with K candidates instead of 1, we could extend ICM and Gibbs sampling to work with more candidates. This leads us to the area of particle swarm optimization, which includes genetic algorithms, which is beyond the scope of this course.

Question

Which of the following algorithms are guaranteed to find the maximum weight assignment (select all that apply)?

backtracking search
greedy search
beam search
Iterated Conditional Modes
Gibbs sampling



Summary so far

Algorithms for max-weight assignments in factor graphs:

(1) Extend partial assignments:

- Backtracking search: exact, exponential time
- Beam search: approximate, linear time

(2) Modify complete assignments:

- Iterated conditional modes: approximate, deterministic
- Gibbs sampling: approximate, randomized



Roadmap

Gibbs Sampling

Conditioning

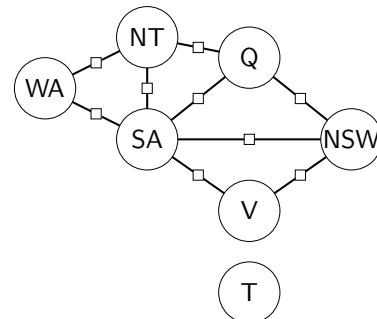
Elimination

Motivation



Key idea: graph

Leverage graph properties to derive efficient algorithms which are exact.



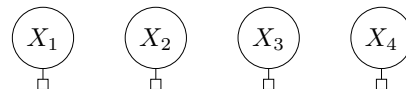
- The goal in the second part of the lecture is to take advantage of the fact that we have a factor **graph**. We will see how exploiting the graph properties can lead us to more efficient algorithms as well as a deeper understanding of the structure of our problem.

- Recall that backtracking search takes time exponential in the number of variables n . While various heuristics can have dramatic speedups in practice, it is not clear how to characterize those improvements rigorously.
- As a motivating example, consider a fully disconnected factor graph. (Imagine n people trying to vote red or blue, but they don't talk to each other.) It's clear that to get the maximum weight assignment, we can just choose the value of each variable that maximizes its own unary factor without worrying about other variables.

Motivation

Backtracking search:

exponential time in number of variables n



Efficient algorithm:

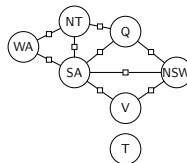
maximize each variable separately

Independence



Definition: independence

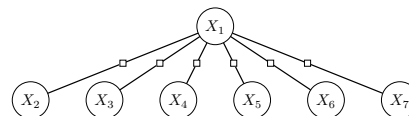
- Let A and B be a partitioning of variables X .
- We say A and B are **independent** if there are no edges between A and B .
- In symbols: $A \perp B$.



$\{WA, NT, SA, Q, NSW, V\}$ and $\{T\}$ are independent.

- Let us formalize this intuition with the notion of **independence**. It turns out that this notion of independence is deeply related to the notion of independence in probability, as we will see in due time.
- Note that we are defining independence purely in terms of the graph structure, which will be important later once we start operating on the graph using two transformations: conditioning and elimination.

Non-independence

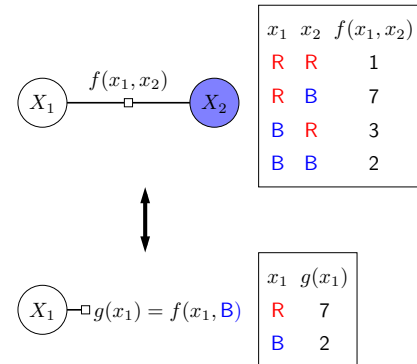


No variables are independent of each other, but feels close...

- When all the variables are independent, finding the maximum weight assignment is easily solvable in time linear in n , the number of variables. However, this is not a very interesting factor graph, because the whole point of a factor graph is to model dependencies (preferences and constraints) between variables.
- Consider the tree-structured factor graph, which corresponds to $n - 1$ people talking only through a leader. Nothing is independent here, but intuitively, this graph should be pretty close to independent.

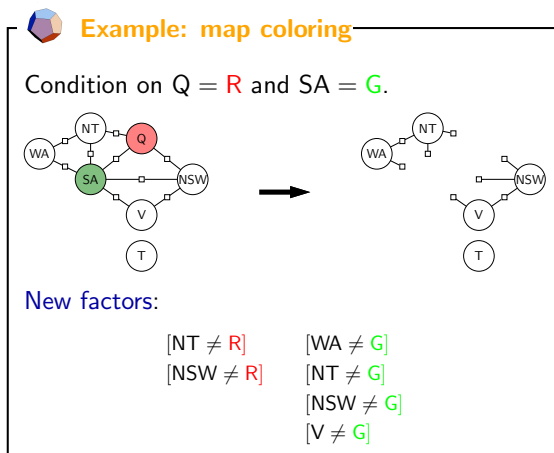
Conditioning

Goal: try to disconnect the graph



Condition on $X_2 = B$: remove X_2, f and add g

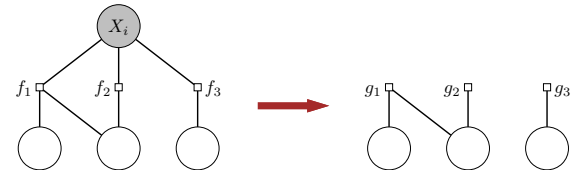
Conditioning: example



- In general, factor graphs are not going to have many partitions which are independent (we got lucky with Tasmania, Australia). But perhaps we can transform the graph to make variables independent. This is the idea of **conditioning**: when we condition on a variable $X_i = v$, this is simply saying that we're just going to clamp the value of X_i to v .
- We can understand conditioning in terms of a graph transformation. For each factor f_j that depends on X_i , we create a new factor g_j . The new factor depends on the scope of f_j excluding X_i ; when called on x , it just invokes f_j with $x \cup \{X_i : v\}$. Think of g_j as a partial evaluation of f_j in functional programming. The transformed factor graph will have each g_j in place of the f_j and also not have X_i .

Conditioning: general

Graphically: remove edges from X_i to dependent factors



Definition: conditioning

- To **condition** on a variable $X_i = v$, consider all factors f_1, \dots, f_k that depend on X_i .
- Remove X_i and f_1, \dots, f_k .
- Add $g_j(x) = f_j(x \cup \{X_i : v\})$ for $j = 1, \dots, k$.

Conditional independence



Definition: conditional independence

- Let A, B, C be a partitioning of the variables.
- We say A and B are **conditionally independent** given C if conditioning on C produces a graph in which A and B are independent.
- In symbols: $A \perp\!\!\!\perp B \mid C$.

Equivalently: every path from A to B goes through C .

Conditional independence

Example: map coloring

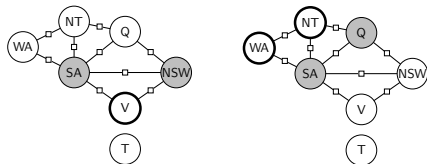
Conditional independence assertion:
 $\{WA, NT\} \perp\!\!\!\perp \{V, NSW, T\} \mid \{SA, Q\}$

CS221 / Summer 2019 / Jia

- With conditioning in hand, we can define **conditional independence**, perhaps the most important property in factor graphs.
- Graphically, if we can find a subset of the variables $C \subset X$ that disconnects the rest of the variables into A and B , then we say that A and B are conditionally independent given C .
- Later, we'll see how this definition relates to the definition of conditional independence in probability.

Markov blanket

How can we separate an arbitrary set of nodes from everything else?

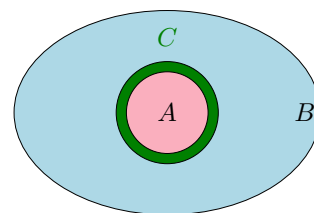


Definition: Markov blanket

Let $A \subseteq X$ be a subset of variables.

Define $\text{MarkovBlanket}(A)$ be the neighbors of A that are not in A .

Markov blanket



Proposition: conditional independence

Let $C = \text{MarkovBlanket}(A)$.

Let $B = X \setminus (A \cup C)$.

Then $A \perp\!\!\!\perp B \mid C$.

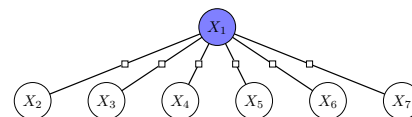
- Suppose we wanted to disconnect a subset of variables $A \subset X$ from the rest of the graph. What is the smallest set of variables C that we need to condition on to make A and the rest of the graph ($B = X \setminus (A \cup C)$) conditionally independent.
- It's intuitive that the answer is simply all the neighbors of A (those that share a common factor) which are not in A . This concept is useful enough that it has a special name: **Markov blanket**.
- Intuitively, the smaller the Markov blanket, the easier the factor graph is to deal with.

Using conditional independence

For each value $v = R, G, B$:

Condition on $X_1 = v$.

Find the maximum weight assignment (easy).



R 3

G 6

B 1

maximum weight is 6

- Now that we understand conditional independence, how is it useful?
- First, this formalizes the fact that if someone tells you the value of a variable, you can condition on that variable, thus potentially breaking down the problem into simpler pieces.
- If we are not told the value of a variable, we can simply try to condition on all possible values of that variable, and solve the remaining problem using any method. If conditioning breaks up the factor graph into small pieces, then solving the problem becomes easier.
- In this example, conditioning on $X_1 = v$ results in a fully disconnected graph, the maximum weight assignment for which can be computed in time linear in the number of variables.

- **Independence** is the key property that allows us to solve subproblems in parallel. It is worth noting that the savings is huge — exponential, not linear. Suppose the factor graph has two disconnected variables, each taking on m values. Then backtracking search would take m^2 time, whereas solving each subproblem separately would take $2m$ time.
- However, the factor graph isn't always disconnected (which would be uninteresting). In these cases, we can **condition** on particular values of a variable. Doing so potentially disconnects the factor graph into pieces, which can be again solved in parallel.
- Factor graphs are interesting because every variable can still influence every other variable, but finding the maximum weight assignment is efficient if there are small bottlenecks that we can condition on.



Summary so far

Independence: when sets of variables A and B are disconnected; can solve separately.

Conditioning: assign variable to value, replaces binary factors with unary factors

Conditional independence: when C blocks paths between A and B

Markov blanket: what to condition on to make A conditionally independent of the rest.

CS221 / Summer 2019 / Jia

37



Roadmap

Gibbs Sampling

Conditioning

Elimination

CS221 / Summer 2019 / Jia

39

Conditioning versus elimination

Conditioning:

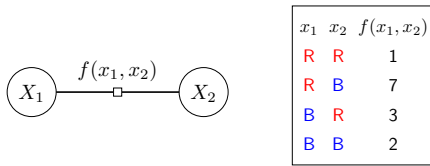
- Removes X_i from the graph
- Add factors that use fixed value of X_i

Elimination (max):

- Removes X_i from the graph
- Add factors that maximize over all values of X_i

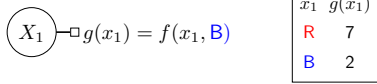
- Now we discuss the second important factor graph transformation: **elimination**. Conditioning was great at breaking the factor graph apart but required a fixed value on which to condition. If we don't know what the value should be, we just have to try all of them.
- Elimination (more precisely, max elimination) also removes variables from the graph, but actually chooses the best value for the eliminated variable X_i . But how do we talk about the best value? The answer is that we compute the best one for all possible assignments to the Markov blanket of X_i . The result of this computation can be stored as a new factor.

Conditioning versus elimination



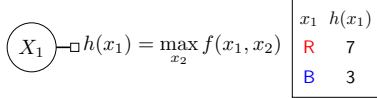
Conditioning:

consider **one** value ($X_2 = B$)



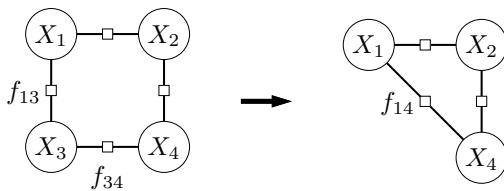
Elimination:

consider **all** values



- If we eliminate X_2 in this simple example, we produce a factor graph with the same structure as what we got for conditioning (but in general, this is not the case), but a different factor.
- In conditioning, the new factor produced $g(x_1)$ was just f evaluated with $x_2 = B$. In elimination, the new factor produced $h(x_1)$ is the maximum value of f over all possible values of x_2 .

Elimination: example



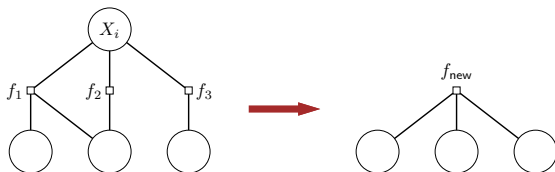
$$f_{14}(x_1, x_4) = \max_{x_3} [f_{13}(x_1, x_3) f_{34}(x_3, x_4)]$$

(maximum weight of assignment to X_3 given X_1, X_4)

$$\max_{x_3} \begin{matrix} x_1 & x_3 & f_{13}(x_1, x_3) \\ \text{R} & \text{R} & 4 \\ \text{R} & \text{B} & 1 \\ \text{B} & \text{R} & 1 \\ \text{B} & \text{B} & 4 \end{matrix} \cdot \begin{matrix} x_3 & x_4 & f_{34}(x_3, x_4) \\ \text{R} & \text{R} & 1 \\ \text{R} & \text{B} & 2 \\ \text{B} & \text{R} & 2 \\ \text{B} & \text{B} & 1 \end{matrix} = \begin{matrix} x_1 & x_4 & f_{14}(x_1, x_4) \\ \text{R} & \text{R} & \max(4 \cdot 1, 1 \cdot 2) = 4 \\ \text{R} & \text{B} & \max(4 \cdot 2, 1 \cdot 1) = 8 \\ \text{B} & \text{R} & \max(1 \cdot 1, 4 \cdot 2) = 8 \\ \text{B} & \text{B} & \max(1 \cdot 2, 4 \cdot 1) = 4 \end{matrix}$$

- Now let us look at a more complex example. Suppose we want to eliminate X_3 . Now we have two factors f_{13} and f_{34} that depend on X_3 .
- Again, recall that we should think of elimination as solving the maximum weight assignment problem over X_3 conditioned on the Markov blanket $\{X_1, X_4\}$.
- The result of this computation is stored in the new factor $f_{14}(x_1, x_4)$, which depends on the Markov blanket.

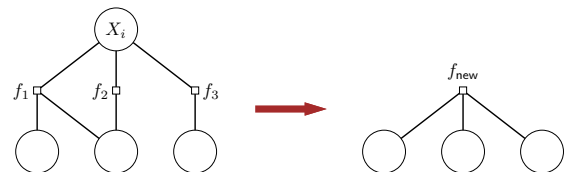
Elimination: general



Definition: elimination

- To **eliminate** a variable X_i , consider all factors f_1, \dots, f_k that depend on X_i .
- Remove X_i and f_1, \dots, f_k .
- Add $f_{\text{new}}(x) = \max_{x_i} \prod_{j=1}^k f_j(x)$

Elimination: general



$$f_{\text{new}}(x) = \max_{x_i} \prod_{j=1}^k f_j(x)$$

- Solves a mini-problem over X_i conditioned on its Markov blanket!
- Scope of f_{new} is $\text{MarkovBlanket}(X_i)$

- In general, to eliminate a variable X_i , we look at all factors which depend on it, just like in conditioning. We then remove those factors f_1, \dots, f_k and X_i , just as in conditioning. Where elimination differs is that it produces a single factor which depends on the Markov blanket rather than a new factor for each f_j .
- Note that eliminating a variable X_i is much more costly than conditioning, and will produce a new factor which can have quite high arity (if X_i depends on many other variables).
- But the good news is that once a variable X_i is eliminated, we don't have to revisit it again. If we have an assignment to the Markov blanket of X_i , then the new factor gives us the weight of the best assignment to X_i , which is stored in the new factor.
- If for every new factor f_{new} , we store for each input, not only the value of the max, but also the argmax, then we can quickly recover the best assignment to X_i .

Question

Suppose we have a star-shaped factor graph. Which of the following is true (select all that apply)?

Conditioning on the hub produces unary factors.

Eliminating the hub produces unary factors.

Variable elimination algorithm



Algorithm: variable elimination

```

For  $i = 1, \dots, n$ :
    Eliminate  $X_i$  (produces new factor  $f_{\text{new},i}$ ).
For  $i = n, \dots, 1$ :
    Set  $X_i$  to the maximizing value in  $f_{\text{new},i}$ .
```

[demo: query(''); maxVariableElimination()]

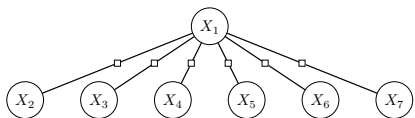
Let max-arity be the maximum arity of any $f_{\text{new},i}$.

Running time: $O(n \cdot |\text{Domain}|^{\text{max-arity}+1})$

- We can turn elimination directly into an actual algorithm for computing the maximum weight assignment by just repeating it until we are left with one variable. This is called the **variable elimination** algorithm.
- The running time of this algorithm is exponential in the maximum arity of the factor produced along the way in variable elimination. The arity in the worst case is $n - 1$, but in the best case it could be a lot better, as we will see.

Variable ordering

What's the maximum arity?



If eliminate leaves first, all factors have arity 1 (good)

If eliminate root first, get giant factor have arity 6 (bad)

Degree heuristic: eliminate variables with the fewest neighbors

- The arity of the factors produced during variable elimination depends on the ordering of the variables. In this extreme example, the difference is between 1 and 6.
- A useful heuristic is to eliminate variables with the smallest Markov blanket. In this example, the heuristic would eliminate the leaves and we'd only end up with factors with arity 1.



Treewidth



Definition: treewidth

The **treewidth** of a factor graph is the maximum arity of any factor created by variable elimination with the **best** variable ordering.

[whiteboard]

- Treewidth of a chain is 1.
- Treewidth of a tree is 1.
- Treewidth of simple cycle is 2.
- Treewidth of $m \times n$ grid is $\min(m, n)$.

- If we use the best ordering, the arity of the largest factor produced is known as the **treewidth**, a very important property in graph theory. Computing the treewidth in general is NP-complete, and verifying that treewidth is k is exponential in k (but linear in the number of nodes).
- However, in practice, it's useful to remember a few examples.
- The treewidth of a chain is 1, by just eliminating all the variables left to right.
- The treewidth of a tree is also 1 by eliminating the variables from the leaves first.
- The treewidth of a simple cycle is 2. By symmetry, we can pick any variable on the cycle; eliminating it creates a factor that connects its two neighbors.
- The treewidth of an $m \times n$ grid is more complex. Without loss of generality, assume that $m \leq n$. One can eliminate the variables by going along the columns left to right and processing the variables from the top row to the bottom row. Verify that when eliminating variable X_{ij} (in the i -th row and the j -th column), its Markov blanket is all the variables in column $j + 1$ and row $\leq i$ as well as all the variables in column j but in row $> i$.
- Note that even if we don't know the exact treewidth, having an upper bound gives us a handle on the running time of variable elimination.



Summary

- **Gibbs sampling**: freely re-assign variables; use randomness to get out of local optima
- **Conditioning**: break up a factor graph into smaller pieces (divide and conquer); can use in backtracking
- **Elimination**: solve a small subproblem conditioned on its Markov blanket

- This lecture, we explored more methods for efficiently finding the maximum weight assignment in a factor graph.
- First, we looked at a second **Local search** method, Gibbs sampling. Like ICM, it modifies the value of one variable at a time. Unlike ICM, it uses randomness to escape local optima and find the global optimum, though it can still get stuck in the worst case.
- The second class of methods are exact methods that rely on (conditional) **independence** structure of the graph, in particular, that the graph is weakly connected, for example, a chain or a tree. We approached this methods by thinking about two graph operations, conditioning and elimination. **Conditioning** sets the value of a variable, and breaks up any factors that touch that variable. **Elimination** maximizes over a variable, but since the maximum value depends on the Markov blanket, this maximization is encoded in a new factor that depends on all the variables in the Markov blanket. The variable elimination computes the maximum weight assignment by applying elimination to each variable sequentially.

Next time

- **First half**: Markov decision processes

State-based models with random transitions

- **Second half**: Bayesian networks

Variable-based models that describe random processes