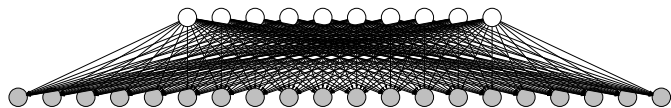




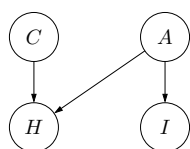
## Lecture 14: Bayesian networks II



## Announcements

- Project progress report due Thursday at **11:00pm**. See Piazza for guidelines.
- **Exam**: Covers up to and including this lecture. Please see Piazza for the logistics.
- Homework 6 (scheduling), Thursday 11:00pm is 2-late-days **hard deadline**.

## Review: Bayesian network



$$\mathbb{P}(C = c, A = a, H = h, I = i) \\ \stackrel{\text{def}}{=} p(c)p(a)p(h | c, a)p(i | a)$$



### Definition: Bayesian network

Let  $X = (X_1, \dots, X_n)$  be random variables.

A **Bayesian network** is a directed acyclic graph (DAG) that specifies a **joint distribution** over  $X$  as a product of **local conditional distributions**, one for each node:

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n p(x_i | x_{\text{Parents}(i)})$$

## Review: probabilistic inference

### Input

Bayesian network:  $\mathbb{P}(X_1 = x_1, \dots, X_n = x_n)$   
 Evidence:  $E = e$  where  $E \subseteq X$  is subset of variables  
 Query:  $Q \subseteq X$  is subset of variables



### Output

$$\mathbb{P}(Q = q | E = e) \text{ for all values } q$$

**Example**: if coughing but no itchy eyes, have a cold?

$$\mathbb{P}(C | H = 1, I = 0)$$

- Last time, we talked about Bayesian networks, which was a fun and convenient modeling framework. We posit a collection of variables that describe the state of the world, and then create a story on how the variables are generated (recall the probabilistic program interpretation).
- Think factor graphs + probability.
- A Bayesian network specifies two parts: (i) a graph structure which governs the qualitative relationship between the variables, and (ii) local conditional distributions, which specify the quantitative relationship.
- Formally, a Bayesian network defines a **joint** probability distribution over many variables (e.g.,  $\mathbb{P}(C, A, H, I)$ ) via the **local** conditional distributions (e.g.,  $p(i | a)$ ). This joint distribution specifies all the information we know about how the world works.
- Think of the joint probability distribution defined by the Bayesian network as a guru. Probabilistic inference allows you to ask the guru anything: what is the probability of having a cold? What if I'm coughing? What if I don't have itchy eyes?
- In the last lecture, we performed probabilistic inference manually. In this lecture, we're going to build a guru that can answer these queries efficiently and automatically.



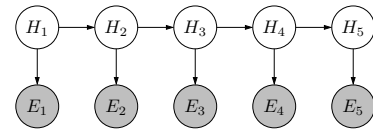
## Roadmap

Forward-backward

Particle filtering

Gibbs sampling

## Hidden Markov model



**Problem: object tracking**

$H_i \in \{1, \dots, K\}$ : location of object at time step  $i$

$E_i \in \{1, \dots, K\}$ : sensor reading at time step  $i$

Start  $p(h_1)$ : e.g., uniform over all locations

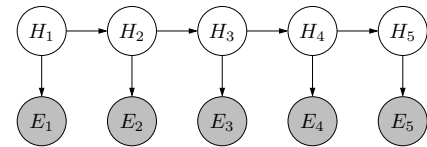
Transition  $p(h_i | h_{i-1})$ : e.g., uniform over adjacent loc.

Emission  $p(e_i | h_i)$ : e.g., uniform over adjacent loc.

$$\mathbb{P}(H = h, E = e) = \underbrace{p(h_1)}_{\text{start}} \prod_{i=2}^n \underbrace{p(h_i | h_{i-1})}_{\text{transition}} \prod_{i=1}^n \underbrace{p(e_i | h_i)}_{\text{emission}}$$

- So far, we have computed various ad-hoc probabilistic queries on ad-hoc Bayesian networks by hand. We will now switch gears and focus on the popular hidden Markov model (HMM), and show how the forward-backward algorithm can be used to compute typical queries of interest.
- As motivation, consider the problem of tracking an object. The probabilistic story is as follows: An object starts at  $H_1$  uniformly drawn over all possible locations. Then at each time step thereafter, it **transitions** to an adjacent location (e.g., with equal probability). For example, if  $H_2 = 3$ , then  $H_3 \in \{2, 4\}$  with equal probability. At each time step, we obtain a sensor reading  $E_i$  given  $H_i$  (e.g., also uniform over locations adjacent to  $H_i$ ).

## Hidden Markov model inference



**Question (filtering):**

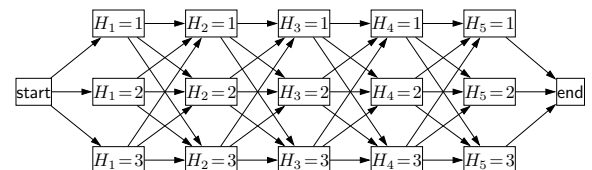
$$\mathbb{P}(H_3 | E_1 = e_1, E_2 = e_2, E_3 = e_3)$$

**Question (smoothing):**

$$\mathbb{P}(H_3 | E_1 = e_1, E_2 = e_2, E_3 = e_3, E_4 = e_4, E_5 = e_5)$$

- In principle, you could ask any type of questions on an HMM, but there are two common ones: filtering and smoothing.
- Filtering asks for the distribution of some hidden variable  $H_i$  conditioned on only the evidence up until that point. This is useful when you're doing real-time object tracking, and you can't see the future.
- Smoothing asks for the distribution of some hidden variable  $H_i$  conditioned on all the evidence, including the future. This is useful when you have collected all the data and want to retroactively go and figure out what  $H_i$  was.

## Lattice representation



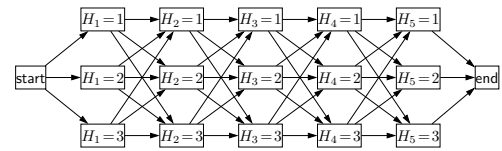
- Edge  $\text{start} \Rightarrow H_1 = h_1$  has weight  $p(h_1)p(e_1 | h_1)$
- Edge  $H_{i-1} = h_{i-1} \Rightarrow H_i = h_i$  has weight  $p(h_i | h_{i-1})p(e_i | h_i)$
- Each path from  $\text{start}$  to  $\text{end}$  is an assignment with weight equal to the product of node/edge weights

- Now let's actually compute these queries. We will do smoothing first. Filtering is a special case: if we're asking for  $H_i$  given  $E_1, \dots, E_i$ , then we can marginalize out the future, reducing the problem to a smaller HMM.
- A useful way to think about inference is returning to state-based models. Consider a graph with a start node, an end node, and a node for each assignment of a value to a variable  $H_i = v$ . The nodes are arranged in a lattice, where each column corresponds to one variable  $H_i$  and each row corresponds to a particular value  $v$ . Each path from the start to the end corresponds exactly to a complete assignment to the nodes.
- Note that in the reduction from a variable-based model to a state-based model, we have committed to an ordering of the variables.
- Each edge has a weight (a single number) determined by the local conditional probabilities (more generally, the factors in a factor graph). For each edge into  $H_i = h_i$ , we multiply by the transition probability into  $h_i$  and emission probability  $p(e_i | h_i)$ . This defines a weight for each path (assignment) in the graph equal to the joint probability  $P(H = h, E = e)$ .
- Note that the lattice contains  $O(nK)$  nodes and  $O(nK^2)$  edges, where  $n$  is the number of variables and  $K$  is the number of values in the domain of each variable.

- The point of bringing back the search-based view is that we can cast the probability queries we care about in terms of sums over paths, and effectively use dynamic programming.
- First, define the forward message  $F_i(v)$  to be the sum of the weights over all paths from the start node to  $H_i = v$ . This can be defined recursively: any path that goes  $H_i = h_i$  will have to go through some  $H_{i-1} = h_{i-1}$ , so we can sum over all possible values of  $h_{i-1}$ .
- Analogously, let the backward message  $B_i(v)$  be the sum of the weights over all paths from  $H_i = v$  to the end node.
- Finally, define  $S_i(v)$  to be the sum of the weights over all paths from the start node to the end node that pass through the intermediate node  $H_i = v$ . This quantity is just the product of the weights of paths going into  $H_i = h_i$  ( $F_i(h_i)$ ) and those leaving it ( $B_i(h_i)$ ).

- Let us go back to the smoothing queries:  $\mathbb{P}(H_i = h_i | E = e)$ . This is just gotten by normalizing  $S_i$ .
- The algorithm is thus as follows: for each node  $H_i = h_i$ , we compute three numbers:  $F_i(h_i)$ ,  $B_i(h_i)$ ,  $S_i(h_i)$ . First, we sweep forward to compute all the  $F_i$ 's recursively. At the same time, we sweep backward to compute all the  $B_i$ 's recursively. Then we compute  $S_i$  by pointwise multiplication.
- Implementation note: we technically can normalize  $S_i$  to get  $\mathbb{P}(H_i | E = e)$  at the very end but it's useful to normalize  $F_i$  and  $B_i$  at each step to avoid underflow. In addition, normalization of the forward messages yields  $\mathbb{P}(H_i = v | E_1 = e_1, \dots, E_i = e_i) \propto F_i(v)$ .

## Lattice representation



**Forward:**  $F_i(h_i) = \sum_{h_{i-1}} F_{i-1}(h_{i-1})w(h_{i-1}, h_i)$

sum of weights of paths from **start** to  **$H_i = h_i$**

**Backward:**  $B_i(h_i) = \sum_{h_{i+1}} B_{i+1}(h_{i+1})w(h_i, h_{i+1})$

sum of weights of paths from  **$H_i = h_i$**  to **end**

**Define**  $S_i(h_i) = F_i(h_i)B_i(h_i)$ :

sum of weights of paths from **start** to **end** through  **$H_i = h_i$**

CS221 / Spring 2020 / Finn & Anari

13

## Lattice representation

**Smoothing queries (marginals):**

$$\mathbb{P}(H_i = h_i | E = e) \propto S_i(h_i)$$



### Algorithm: forward-backward algorithm

Compute  $F_1, F_2, \dots, F_n$

Compute  $B_n, B_{n-1}, \dots, B_1$

Compute  $S_i$  for each  $i$  and normalize

**Running time:**  $O(nK^2)$

CS221 / Spring 2020 / Finn & Anari

15



## Summary

- **Lattice representation:** paths are assignments (think state-based models)
- **Dynamic programming:** compute sums efficiently
- **Forward-backward algorithm:** share intermediate computations across different queries

CS221 / Spring 2020 / Finn & Anari

17



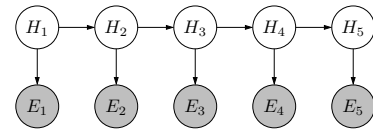
## Roadmap

Forward-backward

**Particle filtering**

Gibbs sampling

## Hidden Markov models



Query (**filtering**):

$$\mathbb{P}(H_1 \mid E_1 = e_1)$$

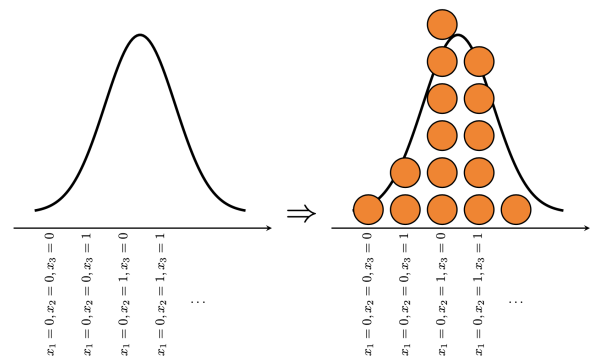
$$\mathbb{P}(H_2 \mid E_1 = e_1, E_2 = e_2)$$

$$\mathbb{P}(H_3 \mid E_1 = e_1, E_2 = e_2, E_3 = e_3)$$

**Motivation:** if  $H_i$  can take on many values, forward-backward is too slow ( $O(nK^2)$ )...

- The forward-backward algorithm runs in  $O(nK^2)$ , where  $K$  is the number of possible values (e.g., locations) that  $H_i$  can take on. This could be a very large number, which makes the forward-backward algorithm very slow (though not exponentially so).
- The motivation of particle filtering is to perform **approximate probabilistic inference**, and leveraging the fact that most of the  $K^2$  pairs are very improbable.
- Although particle filtering applies to general factor graphs, we will present them for hidden Markov models for concreteness.
- As the name suggests, we will use particle filtering for answering filtering queries.

## Approximation by particles



## Review: beam search

Idea: keep  $\leq K$  **candidate list**  $C$  of partial assignments



**Algorithm: beam search**

Initialize  $C \leftarrow \{\emptyset\}$

For each  $i = 1, \dots, n$ :

**Extend:**

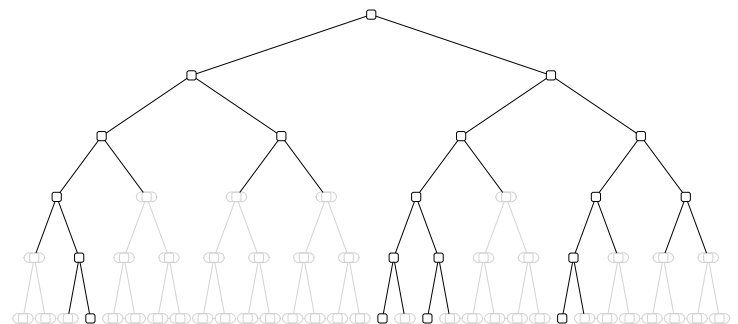
$$C' \leftarrow \{h \cup \{H_i : v\} : h \in C, v \in \text{Domain}_i\}$$

**Prune:**

$$C \leftarrow K \text{ elements of } C' \text{ with highest weights}$$

[demo: beamSearch( $\{K:3\}$ )]

## Review: beam search



Beam size  $K = 4$

- Recall that beam search effectively does a pruned BFS of the search tree of partial assignments, where at each level, we keep track of the  $K$  partial assignments with the highest weight.
- There are two phases. In the first phase, we **extend** all the existing candidates  $C$  to all possible assignments to  $H_i$ ; this results in  $K = |\text{Domain}_i|$  candidates  $C'$ . These  $C'$  are sorted by weight and **pruned** by taking the top  $K$ .

- Beam search does generate a set of particles, but there are two problems.
- First, it can be slow if  $\text{Domain}_i$  is large because we have to try every single value. Perhaps we can be smarter about which values to try.
- Second, we are greedily taking the top  $K$  candidates, which can be too myopic. Can we somehow encourage more diversity?
- Particle filtering addresses both of these problems. There are three steps: propose, which extends the current partial assignment, and weight/resample, which redistributes resources on the particles based on evidence.

- Suppose we have a set of particles that approximates the filtering distribution over  $H_1, H_2$ . The first step is to extend each current partial assignment (particle) from  $h_{1:i-1} = (h_1, \dots, h_{i-1})$  to  $h_{1:i} = (h_1, \dots, h_i)$ .
- To do this, we simply go through each particle and extend it stochastically, using the transition probability  $p(h_i | h_{i-1})$  to sample a new value of  $H_i$ .
- (For concreteness, think of what will happen if  $p(h_i | h_{i-1}) = 0.8$  if  $h_i = h_{i-1}$  and 0.2 otherwise.)
- We can think of advancing each particle according to the dynamics of the HMM. These particles approximate the probability of  $H_1, H_2, H_3$ , but still conditioned on the same evidence.

## Beam search

End result:

- Candidate list  $C$  is set of particles
- Use  $C$  to compute marginals

Problems:

- Extend: slow because requires considering every possible value for  $H_i$
- Prune: greedily taking best  $K$  doesn't provide diversity

Solution (3 steps): **propose, weight, resample**

## Step 1: propose

Old particles:  $\approx \mathbb{P}(H_1, H_2 \mid E_1 = 0, E_2 = 1)$

[0, 1]

[1, 0]



**Key idea: proposal distribution**

For each old particle  $(h_1, h_2)$ , sample  $H_3 \sim p(h_3 \mid h_2)$ .

New particles:  $\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1)$

[0, 1, 1]

[1, 0, 0]

## Step 2: weight

Old particles:  $\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1)$

[0, 1, 1]

[1, 0, 0]



**Key idea: weighting**

For each old particle  $(h_1, h_2, h_3)$ , weight it by  $w(h_1, h_2, h_3) = p(e_3 \mid h_3)$ .

New particles:

$\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1, E_3 = 1)$

[0, 1, 1] (0.8)

[1, 0, 0] (0.4)

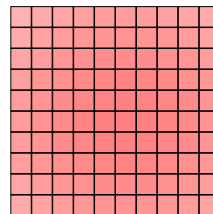
- Having generated a set of  $K$  candidates, we need to now take into account the new evidence  $E_i = e_i$ . This is a deterministic step that simply weights each particle by the probability of generating  $E_i = e_i$ , which is the emission probability  $p(e_i | h_i)$ .
- Intuitively, the proposal was just a guess about where the object will be  $H_3$ . To get a more realistic picture, we condition on  $E_3 = 1$ . Supposing that  $p(e_i = 1 | h_i = 1) = 0.8$  and  $p(e_i = 1 | h_i = 0) = 0.4$ , we then get the weights 0.8 and 0.4. Note that these weights do not sum to 1.

## Step 3: resample

**Question:** given weighted particles, which to choose?

**Tricky situation:**

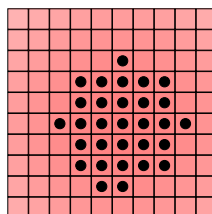
- Target distribution close to uniform
- Fewer particles than locations



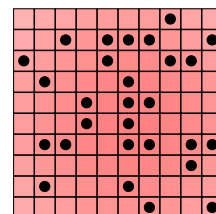
- Having proposed extensions to the particles and computed a weight for each particle, we now come to the question of which particles to keep.
- Intuitively, if a particle has very small weight, then we might want to prune it away. On the other hand, if a particle has high weight, maybe we should dedicate more resources to it.
- As a motivating example, consider an almost uniform distribution over a set of locations, and trying to represent this distribution with fewer particles than locations. This is a tough situation to be in.

## Step 3: resample

$K$  with highest weight



$K$  sampled from distribution



**Intuition:** top  $K$  assignments not representative.

Maybe random samples will be more representative...

- Beam search, which would choose the  $K$  locations with the highest weight, would clump all the particles near the mode. This is risky, because we have no support out farther from the center, where there is actually substantial probability.
- However, if we sample from the distribution which is proportional to the weights, then we can hedge our bets and get a more representative set of particles which cover the space more evenly.

## Step 3: resample



**Key idea: resampling**

Given a distribution  $\mathbb{P}(A = a)$  with  $n$  possible values, draw a sample  $K$  times.

**Intuition:** redistribute particles to more promising areas



**Example: resampling**

$a$	$\mathbb{P}(A = a)$
a1	0.70
a2	0.20
a3	0.05
a4	0.05



sample 1	a1
sample 2	a2
sample 3	a1
sample 4	a1

- After proposing and weighting, we end up with a set of samples  $h_{1:i}$ , each with some weight  $w(h_{1:i})$ . Intuitively, if  $w(h_{1:i})$  is really small, then it might not be worth keeping that particle around.
- Resampling allows us to put (possibly multiple) particles on high weight particles. In the example above, we don't sample a3 and a4 because they have low probability of being sampled.

- In our example, we normalize the particle weights to form a distribution over particles. Then we draw  $K = 2$  independent samples from that distribution. Higher weight particles will get more samples.

## Step 3: resample

Old particles:

$$\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1, E_3 = 1)$$

$$[0, 1, 1] (0.8) \Rightarrow 2/3$$

$$[1, 0, 0] (0.4) \Rightarrow 1/3$$

New particles:

$$\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1, E_3 = 1)$$

$$[0, 1, 1]$$

$$[0, 1, 1]$$

## Particle filtering



**Algorithm: particle filtering**

```

Initialize  $C \leftarrow \{\overbrace{\emptyset, \dots, \emptyset}^{K \text{ times}}\}$ 
For each  $i = 1, \dots, n$ :
  Propose (extend):
     $C' \leftarrow \{h \cup \{H_i : h_i\} : h \in C, h_i \sim p(h_i \mid h_{i-1})\}$ 
  Reweight:
    Compute weights  $w(h) = p(e_i \mid h_i)$  for  $h \in C'$ 
  Resample (prune):
     $C \leftarrow K$  elements drawn independently from  $\propto w(h)$ 

```

[demo: particleFiltering({K:100})]

- The final algorithm here is very similar to beam search. We go through all the variables  $H_1, \dots, H_n$ .
- For each candidate  $h_{i-1} \in C$ , we propose  $h_i$  according to the transition distribution  $p(h_i \mid h_{i-1})$ .
- We then weight this particle using  $w(h) = p(e_i \mid h_i)$ .
- Finally, we select  $K$  particles from  $\propto w(h)$  by sampling  $K$  times independently.

## Particle filtering: implementation

- If only care about last  $H_i$ , collapse all particles with same  $H_i$

$$001 \Rightarrow 1$$

$$101 \Rightarrow 1$$

$$010 \Rightarrow 0$$

$$010 \Rightarrow 0$$

$$110 \Rightarrow 0$$

- If many particles are the same, can just store counts

$$1$$

$$1$$

$$0 \Rightarrow 1 : 2$$

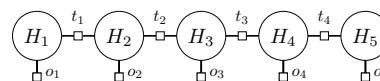
$$0 \Rightarrow 0 : 3$$

$$0$$

$$0$$

- In particle filtering as it is currently defined, each particle is an entire trajectory in the context of object tracking (assignment to all the variables).
- Often in tracking applications, we only care about the last location  $H_i$ , and the HMM is such that the future  $(H_{i+1}, \dots, H_n)$  is conditionally independent of  $H_1, \dots, H_{i-1}$  given  $H_i$ . Therefore, we often just store the value of  $H_i$  rather than its entire ancestry.
- When we only keep track of the  $H_i$ , we might have many particles that have the same value, so it can be useful to store just the counts of each value rather than having duplicates.

## Application: object tracking



### Example: object tracking

- $H_i$ : position of object at  $i$
- Transitions:  $t_i(h_i, h_{i-1}) = [h_i \text{ near } h_{i-1}]$
- Observations:  $o_i(h_i) = \text{sensor reading...}$

$$\mathbb{P}(H = h) \propto \prod_{i=1}^n o_i(h_i) t_i(h_i, h_{i-1})$$

## Particle filtering demo

[see web version]

- Consider a tracking application where an object is moving around in a grid and we are trying to figure out its location  $H_i \in \{1, \dots, \text{grid-width}\} \times \{1, \dots, \text{grid-height}\}$ .
- The transition factors say that from one time step to the next, the object is equally likely to have moved north, south, east, west, or stayed put.
- Each observation is a location on the grid (a yellow dot). The observation factor is a user-defined function which depends on the vertical and horizontal distance.
- Play around with the demo to get a sense of how particle filtering works, especially the different observation factors.



## Roadmap

Forward-backward

Particle filtering

**Gibbs sampling**

## Gibbs sampling

Setup:

Weight( $x$ )



### Algorithm: Gibbs sampling

Initialize  $x$  to a random complete assignment

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

    Compute weight of  $x \cup \{X_i : v\}$  for each  $v$

    Choose  $x \cup \{X_i : v\}$  with probability prop. to weight

[demo]



- Recall that Gibbs sampling proceeds by going through each variable  $X_i$ , considering all the possible assignments of  $X_i$  with some  $v \in \text{Domain}_i$ , computing the weight of the resulting assignment  $x \cup \{X_i : v\}$ , and choosing an assignment with probability proportional to the weight.
- We first introduced Gibbs sampling in the context of local search to get out of local optima.

## Gibbs sampling

Setup:

$$\mathbb{P}(X = x) \propto \text{Weight}(x)$$

### Gibbs sampling (probabilistic interpretation)

Initialize  $x$  to a random complete assignment

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

Set  $X_i = v$  with prob.  $\mathbb{P}(X_i = v \mid X_{-i} = x_{-i})$   
 ( $X_{-i}$  denotes all variables except  $X_i$ )

- Now, we will use it for its original purpose, which is to draw samples from a probability distribution. (In particular, our goal is to draw from the joint distribution over  $X_1, X_2, \dots, X_n$ .) To do this, we need to define a probability distribution given an arbitrary factor graph. Recall that a general factor graph defines a  $\text{Weight}(x)$ , which is the product of all its factors. We can simply normalize the weight to get a distribution:  $\mathbb{P}(X = x) \propto \text{Weight}(x)$ . Then, Gibbs sampling is exactly sampling according to the conditional distribution  $\mathbb{P}(X_i = v \mid X_{-i} = x_{-i})$ .
- Note that the convergence criteria is on the *distribution* of the sample values, not the values themselves.
- Advanced: Gibbs sampling is an instance of a Markov Chain Monte Carlo (MCMC) algorithm which generates a sequence of particles  $X^{(1)}, X^{(2)}, X^{(3)}, \dots$ . A Markov chain is irreducible if there is positive probability of getting from any assignment to any other assignment (now the probabilities are over the random choices of the sampler). When the Gibbs sampler is irreducible, then in the limit as  $t \rightarrow \infty$ , the distribution of  $X^{(t)}$  converges to the true distribution  $\mathbb{P}(X)$ . MCMC is a very rich topic which we will not talk about very much here.

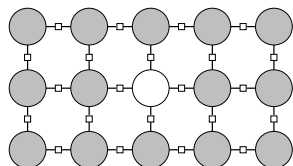
## Application: image denoising



## Application: image denoising



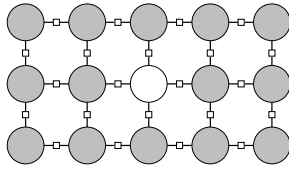
### Example: image denoising



- $X_i \in \{0, 1\}$  is pixel value in location  $i$
- Subset of pixels are observed  
 $o_i(x_i) = [x_i = \text{observed value at } i]$
- Neighboring pixels more likely to be same than different  
 $t_{ij}(x_i, x_j) = [x_i = x_j] + 1$

- Factor graphs (sometimes referred to as Markov random fields) play an important role in computer vision applications. Here we take a look at a very simple image denoising application and construct a factor graph.
- We assume that we have observed some fraction of the pixels in an image, and we wish to recover the pixels which have been removed. Our simple factor graph has two types of factors.
- First,  $o_i(x_i) = 0$  if the pixel at  $i$  is observed and  $x_i$  does not match it. This is analogous to the emission probabilities in an HMM, but here, the factor is deterministic (0 or 1).
- Second,  $t_{ij}(x_i, x_j)$  exists for every pair of neighboring pixels, and encourages them to agree (both be 0 or both be 1). Weight 2 is given to those pairs which are the same and 1 if the pair is different. This is analogous to the transition probabilities in an HMM.

## Application: image denoising



[whiteboard]



### Example: image denoising

If neighbors are 1, 1, 1, 0 and  $X_i$  not observed:

$$\mathbb{P}(X_i = 1 \mid X_{-i} = x_{-i}) = \frac{2 \cdot 2 \cdot 2 \cdot 1}{2 \cdot 2 \cdot 2 \cdot 1 + 1 \cdot 1 \cdot 1 \cdot 2} = 0.8$$

If neighbors are 0, 1, 0, 1 and  $X_i$  not observed:

$$\mathbb{P}(X_i = 1 \mid X_{-i} = x_{-i}) = \frac{1 \cdot 2 \cdot 1 \cdot 2}{1 \cdot 2 \cdot 1 \cdot 2 + 2 \cdot 1 \cdot 1 \cdot 2} = 0.5$$

- Let us compute the Gibbs sampling update. We go through each pixel  $X_i$  and try to update its value. Specifically, we condition on  $X_{-i} = x_{-i}$  being the current value.
- To compute the conditional  $\mathbb{P}(X_i = x_i \mid X_{-i} = x_{-i})$ , we only need to multiply in the factors that are touching  $X_i$ , since conditioning disconnects the graph completely. Many of the factors simply can be ignored.
- For example, suppose we are updating variable  $X_i$ , which has neighbors  $X_j, X_k, X_l, X_m$ . Suppose we condition on  $X_j = 1, X_k = 1, X_l = 0, X_m = 1$ . Then  $\mathbb{P}(X_i = x_i \mid X_{-i} = x_{-i}) \propto t_{ij}(x_i, 1)t_{ik}(x_i, 1)t_{il}(x_i, 0)t_{im}(x_i, 1)o_i(x_i)$ . We can compute the RHS for both  $x_i = 1$  and  $x_i = 0$ , and then normalize to get the conditional probability 0.8.
- Intuitively, the neighbors are all trying to pull  $X_i$  towards their values, and 0.8 reflects the fact that the pull towards 1 is stronger.

## Gibbs sampling: demo

[see web version]

- Try playing with the demo by modifying the settings to get a feeling for what Gibbs sampling is doing. Each iteration corresponds to resampling each pixel (variable).
- When you hit ctrl-enter for the first time, red and black correspond to 1 and 0, and white corresponds to unobserved.
- `showMarginals` allows you to either view the particles produced or the marginals estimated from the particles (this gives you a smoother probability estimate of what the pixel values are).
- If you increase `missingFrac`, the problem becomes harder.
- If you set `coherenceFactor` to 1, this is equivalent to turning off the edge factors.
- If you set `icm` to true, we will use local search rather than Gibbs sampling, which produces very bad solutions.
- In summary, Gibbs sampling is an algorithm that can be applied to any factor graph. It performs a guided random walk in the space of all possible assignments. Under some mild conditions, Gibbs sampling is guaranteed to converge to a sample from the true distribution  $\mathbb{P}(X = x) \propto \text{Weight}(x)$  but this might take until the heat death of the universe for complex models. Nonetheless, Gibbs sampling is widely used due to its simplicity and can work reasonably well.

## Probabilistic inference

Model (Bayesian network or factor graph):

$$\mathbb{P}(X = x) = \prod_{i=1}^n p(x_i \mid x_{\text{Parents}(i)})$$

Probabilistic inference:

$$\mathbb{P}(Q \mid E = e)$$

Algorithms:

- Forward-backward: HMMs, exact
- Particle filtering: HMMs, approximate
- Gibbs sampling: general, approximate

Next time: learning