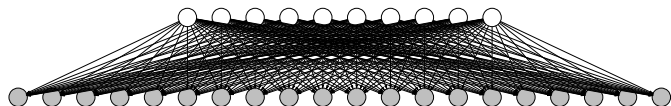




Lecture 7.2: Bayesian networks I



- Situations like these arise all the time in practice: we have a lot of unknowns which are all dependent on one another. If we obtain evidence on some of these unknowns, how does that affect our belief about the other unknowns?
- In this lecture, we'll see how we can perform this type of reasoning under uncertainty in a principled way using Bayesian networks.

- Last time, we talked about factor graphs, which use local factors to specify a weight $\text{Weight}(x)$ for each assignment x in a compact way. The stated objective was to find the maximum weight assignment.
- Given any factor graph, we saw a number of algorithms (backtracking search, beam search, Gibbs sampling, variable elimination) for (approximately) optimizing this objective.

Question

Earthquakes and burglaries are independent events that will cause an alarm to go off. Suppose you hear an alarm. How does hearing on the radio that there's an earthquake change your beliefs?

it increases the probability of burglary

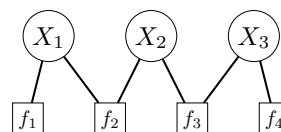
it decreases the probability of burglary

it does not change the probability of burglary

CS221 / Summer 2019 / Jia

1

Review: definition



Definition: factor graph

Variables:

$X = (X_1, \dots, X_n)$, where $X_i \in \text{Domain}_i$

Factors:

f_1, \dots, f_m , with each $f_j(X) \geq 0$

$$\text{Weight}(x) = \prod_{j=1}^m f_j(x)$$

CS221 / Summer 2019 / Jia

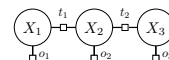
3

Review: person tracking



Problem: person tracking

Sensors report positions: 0, 2, 2. Objects don't move very fast and sensors are a bit noisy. What path did the person take?



- Variables X_i : location of object at time i
- Transition factors $t_i(x_i, x_{i+1})$: incorporate physics
- Observation factors $o_i(x_i)$: incorporate sensors

[demo: `maxVariableElimination()`]

What do the factors **mean**?

CS221 / Summer 2019 / Jia

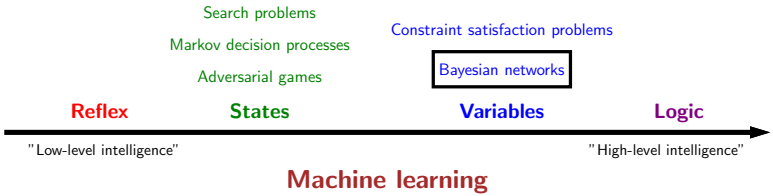
5

- As an example, recall the object tracking example. We defined observation factors to capture the fact that the true object position is close to the sensor reading, and the transition factors to capture the fact that the true object positions across time are close to each other.
- We just set them rather arbitrarily. Is there a more principled way to think about these factors beyond being non-negative functions?

- Much of this class has been on developing modeling frameworks. We started with state-based models, where we cast real-world problems as finding paths or policies through a state graph.
- Then, we saw that for a large class of problems (such as scheduling), it was much more convenient to use the language of factor graphs.
- While factor graphs could be reduced to state-based models by fixing the variable ordering, we saw that they also led to notions of treewidth and variable elimination, which allowed us to understand our models much better.
- In this lecture, we will introduce another modeling framework, Bayesian networks, which are factor graphs imbued with the language of probability. This will give probabilistic life to the factors of factor graphs.

- Bayesian networks were popularized in AI by Judea Pearl in the 1980s, who showed that having a coherent probabilistic framework is important for **reasoning under uncertainty**.
- There is a lot to say about the Bayesian networks (CS228 is an entire course about them and their cousins, Markov networks). So we will devote most of this lecture focusing on modeling.

Course plan



Roadmap

Basics

Probabilistic programs



Review: probability (example)

Random variables: sunshine $S \in \{0, 1\}$, rain $R \in \{0, 1\}$

Joint distribution:

s	r	$\mathbb{P}(S = s, R = r)$
0	0	0.20
0	1	0.08
1	0	0.70
1	1	0.02

Marginal distribution:

s	$\mathbb{P}(S = s)$
0	0.28
1	0.72

(aggregate rows)

Conditional distribution:

s	$\mathbb{P}(S = s \mid R = 1)$
0	0.8
1	0.2

(select rows, normalize)

- Before introducing Bayesian networks, let's review probability (at least the relevant parts). We start with an example about the weather. Suppose we have two boolean random variables, S and R representing sunshine and rain. Think of an assignment to (S, R) as representing a possible state of the world.
- The **joint distribution** specifies a probability for each assignment to (S, R) (state of the world). We use lowercase letters (e.g., s and r) to denote values and uppercase letters (e.g., S and R) to denote random variables. Note that $\mathbb{P}(S = s, R = r)$ is a probability (a number) while $\mathbb{P}(S, R)$ is a distribution (a table of probabilities). We don't know what state of the world we're in, but we know what the probabilities are (there are no unknown unknowns). The joint distribution contains all the information and acts as the central source of truth.
- From it, we can derive a **marginal distribution** over a subset of the variables. We get this by aggregating the rows that share the same value of S . The interpretation is that we are interested in S . We don't explicitly care about R , but we want to take into account R 's effect on S . We say that R is **marginalized out**. This is a special form of elimination. In the last lecture, we leveraged max-elimination, where we took the max over the eliminated variables; here, we are taking a sum.
- The **conditional distribution** selects rows of the table matching the condition (right of the bar), and then normalizes the probabilities so that they sum to 1. The interpretation is that we observe the condition ($R = 1$) and are interested in S . This is the conditioning that we saw for factor graphs, but where we normalize the selected rows to get probabilities.

- In general, we have n random variables X_1, \dots, X_n and let X denote all of them. Suppose X is partitioned into A and B (e.g., $A = (X_1, X_3)$ and $B = (X_2, X_4, X_5)$ if $n = 5$).
- The marginal and conditional distributions can be defined over the subsets A and B rather than just single variables.
- Of course, we can also have a hybrid too: for $n = 3$, $\mathbb{P}(X_1 \mid X_3 = 1)$ marginalizes out X_2 and conditions on $X_3 = 1$.
- It is important to remember the types of objects here: $\mathbb{P}(A)$ is a table where rows are possible assignments to A , whereas $\mathbb{P}(A = a)$ is a number representing the probability of the row corresponding to assignment a .

- At this point, you should have all the definitions to compute any marginal or conditional distribution given access to a joint probability distribution. But what is this really doing and how is this useful?
- We should think about each assignment x as a possible state of the world (it's raining, it's not sunny, there is traffic, it is autumn, etc.). Think of the joint distribution as one giant database that contains full information about how the world works.
- In practice, we'd like to ask questions by querying this probabilistic database. First, we observe some evidence, which effectively fixes some of the variables. Second, we are interested in the distribution of some set of variables which we didn't observe. This forms a query, and the process of answering this query (computing the desired distribution) is called **probabilistic inference**.

Review: probability (general)

Random variables:

$X = (X_1, \dots, X_n)$ partitioned into (A, B)

Joint distribution:

$\mathbb{P}(X) = \mathbb{P}(X_1, \dots, X_n)$

Marginal distribution:

$\mathbb{P}(A) = \sum_b \mathbb{P}(A, B = b)$

Conditional distribution:

$\mathbb{P}(A \mid B = b) \propto \mathbb{P}(A, B = b)$

CS221 / Summer 2019 / Jia

13

Probabilistic inference task

Random variables: unknown quantities in the world

$$X = (S, R, T, A)$$

In words:

- Observe evidence (traffic in autumn): $T = 1, A = 1$
- Interested in query (rain?): R

In symbols:

$$\mathbb{P}(\underbrace{R}_{\text{query}} \mid \underbrace{T = 1, A = 1}_{\text{condition}})$$

$(S \text{ is marginalized out})$

CS221 / Summer 2019 / Jia

15

Challenges

Modeling: How to specify a joint distribution $\mathbb{P}(X_1, \dots, X_n)$ **compactly**?

Bayesian networks (factor graphs for probability distributions)

Inference: How to compute queries $\mathbb{P}(R \mid T = 1, A = 1)$ **efficiently**?

Variable elimination, Gibbs sampling, particle filtering (analogue of algorithms for finding maximum weight assignment)

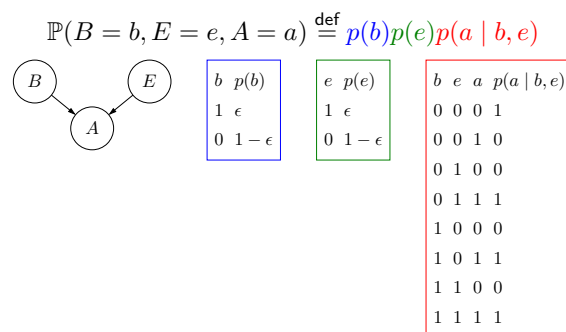
CS221 / Summer 2019 / Jia

17

- In general, a joint distribution over n variables has size exponential in n . From a modeling perspective, how do we even specify an object that large? Here, we will see that Bayesian networks, based on factor graphs, offer an elegant solution.
- From an algorithms perspective, there is still the question of how we perform probabilistic inference efficiently. In the next lecture, we will see how we can adapt all of the algorithms that we saw before for computing maximum weight assignments in factor graphs, essentially by replacing a max with a sum.
- The two desiderata are rather synergistic, and it is the same property — conditional independence — that makes both possible.



Bayesian network (alarm)



$$p(b) = \epsilon \cdot [b = 1] + (1 - \epsilon) \cdot [b = 0]$$

$$p(e) = \epsilon \cdot [e = 1] + (1 - \epsilon) \cdot [e = 0]$$

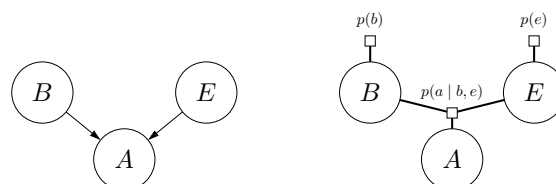
$$p(a | b, e) = [a = (b \vee e)]$$

CS221 / Summer 2019 / Jia

19

- Let us try to model the situation. First, we establish that there are three variables, B (burglary), E (earthquake), and A (alarm). Next, we connect up the variables to model the dependencies.
- Unlike in factor graphs, these dependencies are represented as **directed** edges. You can intuitively think about the directionality as suggesting causality, though what this actually means is a deeper question and beyond the scope of this class.
- For each variable, we specify a **local conditional distribution** (a factor) of that variable given its parent variables. In this example, B and E have no parents while A has two parents, B and E . This local conditional distribution is what governs how a variable is generated.
- We are writing the local conditional distributions using p , while \mathbb{P} is reserved for the joint distribution over all random variables, which is defined as the product.

Bayesian network (alarm)



$$\mathbb{P}(B = b, E = e, A = a) = p(b)p(e)p(a | b, e)$$

Bayesian networks are a special case of factor graphs!

CS221 / Summer 2019 / Jia

21

- Note that the local conditional distributions (e.g., $p(a | b, e)$) are non-negative so they can be thought of simply as factors of a factor graph. The joint probability of an assignment is then the weight of that assignment.
- In this light, Bayesian networks are just a type of factor graphs, but with additional structure and interpretation.

Probabilistic inference (alarm)

Joint distribution:

b	e	a	$\mathbb{P}(B = b, E = e, A = a)$
0	0	0	$(1 - \epsilon)^2$
0	0	1	0
0	1	0	0
0	1	1	$(1 - \epsilon)\epsilon$
1	0	0	0
1	0	1	$\epsilon(1 - \epsilon)$
1	1	0	0
1	1	1	ϵ^2

Queries: $\mathbb{P}(B)$? $\mathbb{P}(B | A = 1)$? $\mathbb{P}(B | A = 1, E = 1)$?

[demo: $\epsilon = 0.05$]

CS221 / Summer 2019 / Jia

23

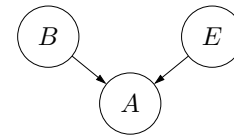
- Bayesian networks can be used to capture common reasoning patterns under uncertainty (which was one of their first applications).
- Consider the following model: Suppose the probability of an earthquake is ϵ and the probability of a burglary is ϵ and both are independent. Suppose that the alarm always goes off if either an earthquake or a burglary occurs.
- In the prior, we can eliminate A and E and get that the probability of the burglary is ϵ .
- Now suppose we hear the alarm $A = 1$. The probability of burglary is now $\mathbb{P}(B = 1 \mid A = 1) = \frac{1}{2-\epsilon}$.
- Now suppose that you hear on the radio that there was an earthquake ($E = 1$). Then the probability of burglary goes down to $\mathbb{P}(B = 1 \mid A = 1, E = 1) = \epsilon$ again.

- This last phenomenon has a special name: **explaining away**. Suppose we have two **cause** variables B and E , which are parents of an **effect** variable A . Assume the causes influence the effect positively (e.g., through the OR function).
- Conditioned on the effect $A = 1$, there is some posterior probability of B . Conditioned on the effect $A = 1$ and the other cause $E = 1$, the new posterior probability is reduced. We then say that the other cause E has explained away B .

- Without further ado, let's define a Bayesian network formally. A Bayesian network defines a large joint distribution in a modular way, one variable at a time.
- First, the graph structure captures what other variables a given variable depends on.
- Second, we specify a local conditional distribution for variable X_i , which is a function that specifies a distribution over X_i given an assignment $x_{\text{Parents}(i)}$ to its parents in the graph (possibly none). The joint distribution is simply **defined** to be the product of all of the local conditional distributions together.
- Notationally, we use lowercase p (in $p(x_i \mid x_{\text{Parents}(i)})$) to denote a local conditional distribution, and uppercase \mathbb{P} to denote the induced joint distribution over all variables. While the two can coincide, it is important to keep these things separate in your head!



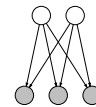
Explaining away



Key idea: explaining away

Suppose two causes positively influence an effect. Conditioned on the effect, conditioning on one cause reduces the probability of the other cause.

Definition



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:

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

Special properties

Key difference from general factor graphs:



Key idea: locally normalized

All factors (local conditional distributions) satisfy:

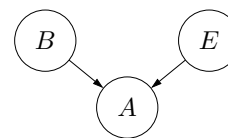
$$\sum_{x_i} p(x_i \mid x_{\text{Parents}(i)}) = 1 \text{ for each } x_{\text{Parents}(i)}$$

Implications:

- Consistency of sub-Bayesian networks
- Consistency of conditional distributions

- But Bayesian networks are more than that. The key property is that all the local conditional distributions, being distributions, sum to 1 over the first argument.
- This simple property results in two important properties of Bayesian networks that are not present in general factor graphs.

Consistency of sub-Bayesian networks



A short calculation:

$$\begin{aligned}
 \mathbb{P}(B = b, E = e) &= \sum_a \mathbb{P}(B = b, E = e, A = a) \\
 &= \sum_a p(b)p(e)p(a | b, e) \\
 &= p(b)p(e) \sum_a p(a | b, e) \\
 &= p(b)p(e)
 \end{aligned}$$

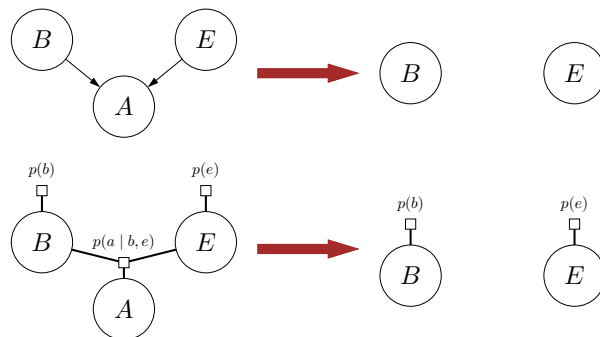
- First, let's see what happens when we marginalize A (by performing algebra on the joint probability). We see that we end up with $p(b)p(e)$, which actually defines a sub-Bayesian network with one fewer variable, and the same local conditional probabilities.
- If one marginalizes out all the variables, then one gets 1, which verifies that a Bayesian network actually defines a probability distribution.
- The philosophical ramification of this property is that there could be many other variables that depend on the variables you've modeled (earthquakes also impacts traffic) but as long as you don't observe them, they can be ignored mathematically (ignorance is bliss). Note that this doesn't mean that knowing about the other things isn't useful.

Consistency of sub-Bayesian networks



Key idea: marginalization

Marginalization of a leaf node yields a Bayesian network without the node.



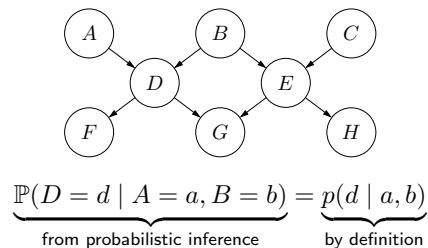
- This property is very attractive, because it means that whenever we have a large Bayesian network, where we don't care about some of the variables, we can just remove them (graph operations), and this encodes the same distribution as we would have gotten from marginalizing out variables (algebraic operations). The former, being visual, can be more intuitive.

Consistency of local conditionals



Key idea: local conditional distributions

Local conditional distributions (factors) are the true conditional distributions.



- Note that the local conditional distributions $p(d \mid a, b)$ are simply defined by the user. On the other hand, the quantity $\mathbb{P}(D = d \mid A = a, B = b)$ is not defined, but follows from probabilistic inference on the joint distribution defined by the Bayesian network.
- It's not clear a priori that the two have anything to do with each other. The second special property that we get from using Bayesian networks is that the two are actually the same.
- To show this, we can remove all the descendants of D by the consistency of sub-Bayesian networks, leaving us with the Bayesian network $\mathbb{P}(A = a, B = b, D = d) = p(a)p(b)p(d \mid a, b)$. By the chain rule, $\mathbb{P}(A = a, B = b, D = d) = \mathbb{P}(A = a, B = b)\mathbb{P}(D = d \mid A = a, B = b)$. If we marginalize out D , then we are left with the Bayesian network $\mathbb{P}(A = a, B = b) = p(a)p(b)$. From this, we can conclude that $\mathbb{P}(D = d \mid A = a, B = b) = p(d \mid a, b)$.
- This argument generalizes to any Bayesian network and local conditional distribution.



Medical diagnosis



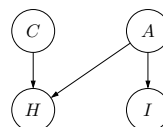
Problem: cold or allergies?

You are coughing and have itchy eyes. Do you have a cold or allergies?

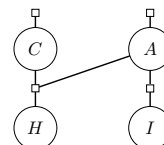
[demo]

Variables: Cold, Allergies, Cough, Itchy eyes

Bayesian network:



Factor graph:



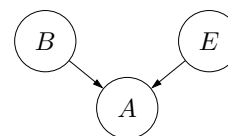
CS221 / Summer 2019 / Jia

37

- Here is another example (a cartoon version of Bayesian networks for medical diagnosis). Allergies and cold are the two hidden variables that we'd like to infer (we have some prior over these two). Cough and itchy eyes are symptoms that we observe as evidence, and we have some likelihood model of these symptoms given the hidden causes.
- We can use the demo to infer the hidden state given the evidence.



Summary so far



- Set of random variables capture state of world
- Local conditional distributions \Rightarrow joint distribution
- Probabilistic inference task: ask questions
- Captures reasoning patterns (e.g., explaining away)
- Factor graph interpretation (for inference later)

CS221 / Summer 2019 / Jia

39

Roadmap

Basics

Probabilistic programs

Probabilistic programs

Goal: make it easier to write down complex Bayesian networks

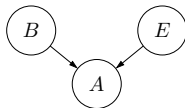


Key idea: probabilistic program

Write a program to generate an assignment (rather than specifying the probability of an assignment).



Probabilistic programs



Probabilistic program: alarm

```

B ~ Bernoulli(ε)
E ~ Bernoulli(ε)
A = B ∨ E

```



Key idea: probabilistic program

A randomized program that sets the random variables.

```

def Bernoulli(epsilon):
    return random.random() < epsilon

```

- There is another way of writing down Bayesian networks other than graphically or mathematically, and that is as a probabilistic program. A **probabilistic program** is a randomized program that invokes a random number generator to make random choices. Executing this program will assign values to a collection of random variables X_1, \dots, X_n ; that is, generating an assignment.
- The probability (e.g., fraction of times) that the program generates that assignment is exactly the probability under the joint distribution specified by that program.
- We should think of this program as outputting the state of the world (or at least the part of the world that we care about for our task).
- Note that the probabilistic program is only used to define joint distributions. We usually wouldn't actually run this program directly.
- For example, we show the probabilistic program for alarm. $B \sim \text{Bernoulli}(\epsilon)$ simply means that $\mathbb{P}(B = 1) = \epsilon$. Here, we can think about $\text{Bernoulli}(\epsilon)$ as a randomized function (`random() < epsilon`) that returns 1 with probability ϵ and 0 with probability $1 - \epsilon$.

Probabilistic program: example



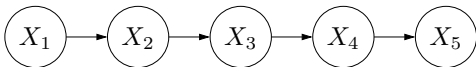
Probabilistic program: object tracking

```

X0 = (0, 0)
For each time step i = 1, ..., n:
    With probability α:
        Xi = Xi-1 + (1, 0) [go right]
    With probability 1 - α:
        Xi = Xi-1 + (0, 1) [go down]

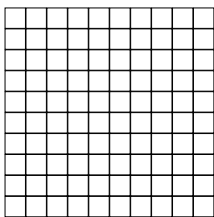
```

Bayesian network structure:



- This is a more interesting generative model since it has a for loop, which allows us to determine the distribution over a templated set of n variables rather than just 3 or 4.
- In these cases, variables are generally indexed by something like time or location.
- We can also draw the Bayesian network. Each X_i only depends on X_{i-1} . This is a chain-structured Bayesian network, called a **Markov model**.

Probabilistic program: example



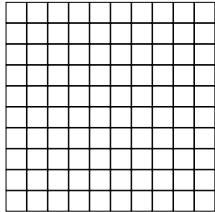
(press ctrl-enter to save)

Run

- Try clicking [Run]. Each time a new assignment of (X_1, \dots, X_n) is chosen.

Probabilistic inference: example

Query: what are possible trajectories given **evidence** $X_{10} = (8, 2)$?



(press ctrl-enter to save)

Run

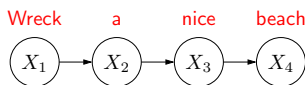
- This program only serves for defining the distribution. Now we can query that distribution and ask the question: suppose the program set $X_{10} = (8, 2)$; what is the distribution over the other variables?
- In the demo, note that all trajectories are constrained to go through $(8, 2)$ at time step 10.

Application: language modeling



Probabilistic program: Markov model

For each position $i = 1, 2, \dots, n$:
Generate word $X_i \sim p(X_i | X_{i-1})$

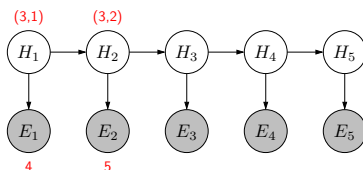


Application: object tracking



Probabilistic program: hidden Markov model (HMM)

For each time step $t = 1, \dots, T$:
Generate object location $H_t \sim p(H_t | H_{t-1})$
Generate sensor reading $E_t \sim p(E_t | H_t)$



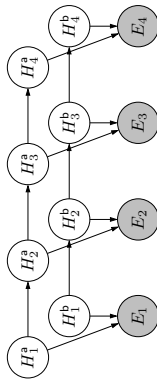
- Markov models are limiting because they do not have a way of talking about noisy evidence (sensor readings). They can be extended quite easily to hidden Markov models, which introduce a parallel sequence of observation variables.
- For example, in object tracking, H_t denotes the true object location, and E_t denotes the noisy sensor reading, which might be (i) the location H_t plus noise, or (ii) the distance from H_t plus noise, depending on the type of sensor.
- In speech recognition, H_t would be the phonemes or words and E_t would be the raw acoustic signal.

Applications: speech recognition, information extraction, gene finding

Application: multiple object tracking

Probabilistic program: factorial HMM

For each time step $t = 1, \dots, T$:
 For each object $o \in \{a, b\}$:
 Generate location $H_t^o \sim p(H_t^o | H_{t-1}^o)$
 Generate sensor reading $E_t \sim p(E_t | H_t^a, H_t^b)$



CS221 / Summer 2019 / Jia

54

- An extension of an HMM, called a **factorial HMM**, can be used to track multiple objects. We assume that each object moves independently according to a Markov model, but that we get one sensor reading which is some noisy aggregated function of the true positions.
- For example, E_t could be the set $\{H_t^a, H_t^b\}$, which reveals where the objects are, but doesn't say which object is responsible for which element in the set.

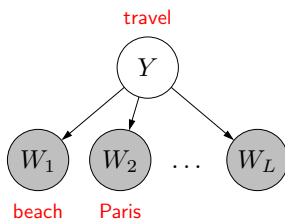
Application: document classification

Question: given a text document, what is it about?



Probabilistic program: naive Bayes

Generate label $Y \sim p(Y)$
 For each position $i = 1, \dots, L$:
 Generate word $W_i \sim p(W_i | Y)$



CS221 / Summer 2019 / Jia

56

- Naive Bayes is a very simple model which can be used for classification. For document classification, we generate a label and all the words in the document given that label.
- Note that the words are all generated independently, which is not a very realistic model of language, but naive Bayes models are surprisingly effective for tasks such as document classification.

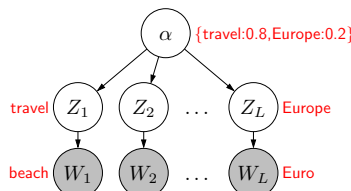
Application: topic modeling

Question: given a text document, what topics is it about?



Probabilistic program: latent Dirichlet allocation

Generate a distribution over topics $\alpha \in \mathbb{R}^K$
 For each position $i = 1, \dots, L$:
 Generate a topic $Z_i \sim p(Z_i | \alpha)$
 Generate a word $W_i \sim p(W_i | Z_i)$



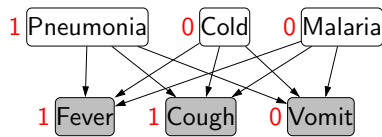
CS221 / Summer 2019 / Jia

58

- A more sophisticated model of text is latent Dirichlet Allocation (LDA), which allows a document to not just be about one topic (which was true in naive Bayes), but about multiple topics.
- Here, the distribution over topics α is chosen per document from a Dirichlet distribution. Note that α is a continuous-valued random variable. For each position, we choose a topic according to that per-document distribution and generate a word given that topic.
- Latent Dirichlet Allocation (LDA) has been very influential for modeling not only text but images, videos, music, etc.; any sort of data with hidden structure. It is very related to matrix factorization.

Application: medical diagnostics

Question: If patient has a cough and fever, what disease(s) does he/she have?



Probabilistic program: diseases and symptoms

For each disease $i = 1, \dots, m$:

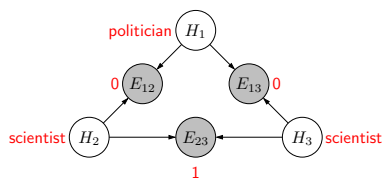
Generate activity of disease $D_i \sim p(D_i)$

For each symptom $j = 1, \dots, n$:

Generate activity of symptom $S_j \sim p(S_j \mid D_{1:m})$

Application: social network analysis

Question: Given a social network (graph over n people), what types of people are there?



Probabilistic program: stochastic block model

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

Generate person type $H_i \sim p(H_i)$

For each pair of people $i \neq j$:

Generate connectedness $E_{ij} \sim p(E_{ij} \mid H_i, H_j)$

- We already saw a special case of this model. In general, we would like to diagnose many diseases and might have measured many symptoms and vitals.

- One can also model graphs such as social networks. A very naive-Bayes-like model is that each node (person) has a "type". Whether two people interact with each other is determined solely by their types and random chance.
- Note: there are extensions called mixed membership models which, like LDA, allow each person to have multiple types.
- In summary, it is quite easy to come up with probabilistic programs that tell a story of how the world works for the domain of interest. These probabilistic programs define joint distributions over assignments to a collection of variables. Usually, these programs describe how some collection of hidden variables H that you're interested in behave, and then describe the generation of the evidence E that you see conditioned on H . After defining the model, one can do probabilistic inference to compute $\mathbb{P}(H \mid E = e)$.



Summary

Bayesian networks: modular definition of large joint distribution over variables

Probabilistic programs: generative stories of how the world works

Next time: probabilistic inference