

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
Department of Electrical Engineering and Computer Science
6.081—Introduction to EECS I
Spring Semester, 2007

Lecture 12 Notes

Discrete Probability and State Estimation

Where am I?

Last week, we introduced the idea of a state space, and its use for planning trajectories from some starting state to a goal. Our assumptions in that work were that we knew the initial state, and that the actions could be executed without error. That is a useful idealization in many cases, but it is also very frequently false. Even navigation through a city can fail on both counts: sometimes we don't know where we are on a map, and sometimes, due to traffic or road work or bad driving, we fail to execute a turn we had intended to take.

In such situations, we have some information about where we are: we can make observations of our local surroundings, which give us useful information; and we know what actions we have taken and the consequences those are likely to have on our location. So, the question is: how can we take information from a sequence of actions and local observations and integrate it into some sort of estimate of where we are? What form should that estimate take?

We'll consider a probabilistic approach to answering this question. We'll assume that, as you navigate, you maintain a *belief state* which contains your best information about what state you're in, which is represented as a probability distribution over all possible states. So, it might say that you're sure you're somewhere in Boston, and you're pretty sure it's Storrow drive, but you don't know whether you're past the Mass Ave bridge or not (of course, it will specify this all much more precisely).

We'll start with some basic background on probability, and then talk about how to estimate an underlying hidden state of a changing world, based on noisy and partial observations.

Probability

Probability theory is a calculus that allows us to assign numerical assessments of uncertainty to possible events, and then do calculations with them in a way that preserves their meaning. (A similar system that you might be more familiar with is algebra: you start with some facts that you know, and the axioms of algebra allow you to make certain manipulations to your equations that you know will preserve their truth).

The typical informal interpretation of probability statements is that they are long-term frequencies: to say “the probability that this coin will come up heads when flipped is 0.5” is to say that, in the long run, the proportion of flips that come up heads will be 0.5. This is known as the *frequentist interpretation* of probability. But then, what does it mean to say “there is a 0.7 probability that it will rain somewhere in Boston sometime on April 29, 2007”? How can we repeat that process a lot of times, when there will only be one April 29, 2007? Another way to interpret probabilities is that

they are measures of a person's (or robot's) *degree of belief* in the statement. This is sometimes referred to as the *Bayesian interpretation*. In the Bayesian interpretation, you cannot be wrong about your beliefs, but it is possible (and suboptimal, from the perspective of maximizing utility) to be inconsistent.

So, studying and applying the axioms of probability will help us make true statements about long-run frequencies and make consistent statements about our beliefs, by deriving sensible consequences from initial assumptions.

We'll just consider the case of discrete sample spaces, so we'll let \mathcal{U} be the *universe* or sample space, which is a set of *atomic events*. An atomic event is just an outcome or a way the world could be. It might be a die roll, or whether the robot is in a particular room, for example. An *event* is a subset of \mathcal{U} . A probability measure P is a mapping from events to numbers that satisfy the following axioms:

$$\begin{aligned} P(\mathcal{U}) &= 1 \\ P(\{\}) &= 0 \\ P(A \cup B) &= P(A) + P(B) - P(A \cap B) \end{aligned}$$

Or, in English:

- The probability that something will happen is 1.
- The probability that nothing will happen is 0.
- The probability that an atomic event in the set A or an atomic event in the set B will happen is the probability that an atomic event of A will happen plus the probability that an atomic event of B will happen, minus the probability that an atomic event that is in both A and B will happen (because those events effectively got counted twice in the sum of $P(A)$ and $P(B)$).

Armed with these axioms, we are prepared to do anything that can be done with discrete probability!

Random variables A discrete random variable is a discrete set of values, $v_1 \dots v_n$ and a mapping of those values to probabilities $p_1 \dots p_n$ such that $p_i \in [0, 1]$ and $\sum_i p_i = 1$. So, for instance, the random variable associated with flipping a somewhat biased coin might be {heads : 0.6, tails : 0.4}.

In a world that is appropriately described with multiple random variables, the atomic event space is the *cartesian product* of the value spaces of the variables. So, for example, consider two random variables, C for *cavity* and A for *toothache*. If they can each take on the values T or F, then the universe is

$$C \times A = \{(T, T), (T, F), (F, T), (F, F)\} .$$

The *joint distribution* of a set of random variables is a function from elements of the product space to probability values that sum to 1 over the whole space. So, for example, consider the following table:

		C		
		T	F	
A	T	0.05	0.05	0.1
	F	0.1	0.8	0.9
		0.15	0.85	

The bold entries in the table make up the joint probability distribution. They are assignments of probability values to atomic events, which are complete specifications of the values of all of the

random variables. For example, $P(C = T, A = F) = 0.1$; that is, the probability of the atomic event that random variable C has value T and random variable A has value F is 0.1. Other events can be made up of the union of these primitive events, and specified by the assignments of values to only some of the variables. So, for instance, the event $A = T$ is really a set of primitive events: $\{(A = T, C = F), (A = T, C = T)\}$, which means that

$$P(A = T) = P(A = T, C = T) + P(A = T, C = F) ,$$

which is just the sum of the row in the table.

The rightmost column of numbers and the bottommost row of numbers are the *marginal probability distributions* of the individual random variables. So, for example, we can see from the marginals that $P(A = T) = 0.1$, and that $P(C = T) = 0.15$. Although you can compute the marginal distributions from the joint distribution, **you cannot in general compute the joint distribution from the marginal distribution!!**.

In the **very special case** when two random variables A and B do not influence one another, we say that they are *independent*, which is mathematically defined as

$$P(A = a, B = b) = P(A = a)P(B = b) .$$

If we only knew the marginals of toothaches and cavities, and assumed they were independent, we would find that $P(C = T, A = T) = 0.015$, which is much less than the value in our table. This is because, although cavity and toothache are relatively rare events, they are highly dependent.

One more important idea is *conditional probability*, where we ask the probability of some event E_1 , assuming that some other event E_2 is true; we do this by restricting our attention to the part of the sample space in E_2 . The conditional probability is the amount of the sample space that is both in E_1 and E_2 , divided by the amount in E_2 :

$$P(E_1|E_2) = \frac{P(E_1, E_2)}{P(E_2)} .$$

So, if a patient walks into your dental practice saying that she has a toothache, what's the probability she has a cavity? That is the conditional probability of $C = T$ given $A = T$ (we already know the value of A , so our only uncertainty is about C).

$$\begin{aligned} P(C = T|A = T) &= \frac{P(C = T, A = T)}{P(A = T)} \\ &= \frac{0.05}{0.1} \\ &= 0.5 \end{aligned}$$

So, although a cavity is relatively unlikely, it becomes much more likely conditioned on knowing that the person has a toothache.

Bayes' Rule Sometimes, for medical diagnosis or characterizing the quality of a sensor, it's easiest to measure conditional probabilities of the form $P(\text{symptom} = T | \text{disease} = T)$, indicating in what proportion of diseased patients a particular symptom shows up. (These numbers are often more useful, because they tend to be the same everywhere, even though the proportion of the population that has disease may differ.) But in these cases, we really want to know $P(\text{disease} =$

$T|_{\text{symptom} = T}$). We can use the definition of conditional probability in a form that is known as *Bayes' Rule* to get this:

$$P(\text{disease} = T | \text{symptom} = T) = \frac{P(\text{symptom} = T | \text{disease} = T)P(\text{disease} = T)}{P(\text{symptom} = T)} .$$

This rule is often very useful, and can easily be verified:

$$\begin{aligned} P(E_1|E_2) &= \frac{P(E_2|E_1)P(E_1)}{P(E_2)} \\ \frac{P(E_1, E_2)}{P(E_2)} &= \frac{P(E_2, E_1)P(E_1)}{P(E_1)P(E_2)} \\ &= \frac{P(E_2, E_1)}{P(E_2)} \end{aligned}$$

State estimation

So, now, let's consider the application of interest: there is a system moving through some sequence of states over time, but instead of getting to see the states, we only get to make a sequence of observations of the system. The question is: what can we infer about the current state of the system given the history of observations we have made?

As a very simple example, let's consider a copy machine: we'll model it in terms of two possible internal states: *good* and *bad*. But since we don't get to see inside the machine, we can only make observations of the copies it generates; they can either be *perfect*, *smudged*, or *all black*.

We can model this problem as a *hidden Markov model* (HMM). In an HMM we model time a discrete sequence of steps, and we can think about the state and the observation at each step. So, we'll use random variables S_0, S_1, S_2, \dots to model the state at each time step and random variables O_1, O_2, \dots to model the observation at each time step. Our problem will be to compute the state at some current time t given the past history of observations; that is, to compute

$$P(S_t | O_1 \dots O_t) .$$

A hidden Markov model makes a strong assumption about the system: that the state at time t is sufficient to determine the probability distribution over the observation at time t and the state at time $t + 1$. Furthermore, we assume that the way the state at time 3 depends on the state at time 2 is the same way that the state at time 2 depends on the state at time 1, and so on; similarly for the observations.

So, in order to specify our model of how this system works, we need to provide three sets of probability distributions:

Initial state distribution: We need to have some idea of the state that the machine will be in at the very first step of time that we're modeling. This is often also called the *prior* distribution on the state. We'll write it as

$$P(S_0 = s) .$$

It will be a collection of probability values, one for each possible state s , that sum to 1.

Initial state distribution			State transition model		
	s			s	
	good	bad		good	bad
$P(S_0 = s)$	0.9	0.1	$P(S_{t+1} = s S_t = \text{good})$	0.7	0.3
			$P(S_{t+1} = s S_t = \text{bad})$	0.1	0.9

Observation model			
	o		
	perfect	smudge	black
$P(O_t = o S_t = \text{good})$	0.8	0.1	0.1
$P(O_t = o S_t = \text{bad})$	0.1	0.7	0.2

Figure 1: An HMM model of the copy-machine diagnosis problem.

State transition model: Next, we need to specify how the state of the system will change over time. We do that by considering each possible state, s_t , that the system could be in at time t , and then writing down the conditional probability distribution over S_{t+1} ,

$$P(S_{t+1} = s_{t+1} | S_t = s_t) ,$$

which specifies for each possible new state, s_{t+1} , how likely it will be, given that the system was in state s_t on the time step before.

Observation model: Finally, we need to specify how the observations we make of the system depend on the underlying state. This is often also called the *sensor model*. We specify it by considering each possible value, s_t , that the system could be in at time t , and then writing down the conditional probability distribution

$$P(O_t = o_t | S_t = s_t) ,$$

which specifies for each possible observation o_t , how likely it will be, given that the system is currently in state s_t .

Copy machine example

Figure 1 shows an HMM model for the copy-machine diagnosis problem. Note that the state transition model and the observation model consist of a set of conditional distributions; each one is represented by a separate row, conditioned on the current state. The numbers in each distribution must, as always, add up to 1.

Our first copy So, now, let's assume we get a brand new copy machine in the mail, and we think it is probably (0.9) good, but we're not entirely sure. We print out a page, and it looks

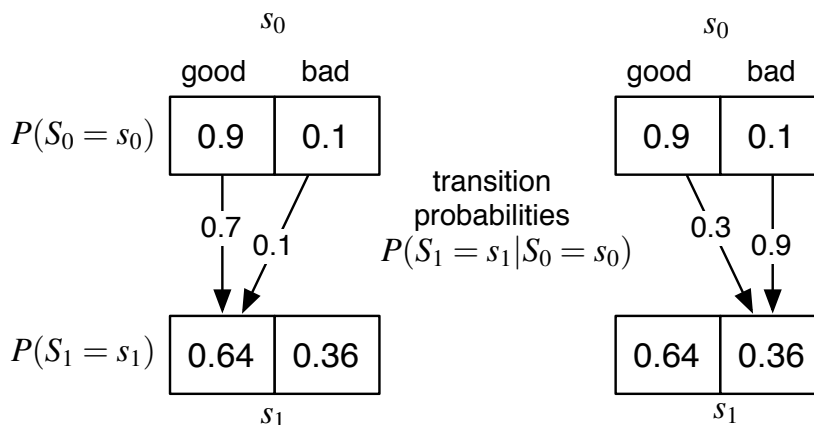


Figure 2: Schematic version of first transition update

perfect. Yay! Now, what do we believe about the state of the machine? We'd like to compute $P(S_1 = \text{good} | O_1 = \text{perfect})$. We'll do this in two steps. First we'll consider how the machine might have changed from time step 0 to 1, and then we'll consider what information we have gained from the observation.

So, we'll start by thinking about $P(S_1 = \text{good})$, imagining that we haven't yet seen the first printout (or that some annoying office-mate took it by mistake). We'll start by realizing that the machine could be in a good state now either because it was in a good state before and stayed good, or because it was in a bad state before, and magically repaired itself (ha!):

$$P(S_1 = \text{good}) = P(S_1 = \text{good}, S_0 = \text{good}) + P(S_1 = \text{good}, S_0 = \text{bad}) .$$

Now, our model doesn't give us the probability of the state at one step *and* the state at the next step; but we can use the definition of conditional probability to rewrite $P(E_1, E_2) = P(E_1 | E_2)P(E_2)$, which lets us write

$$P(S_1 = \text{good}) = P(S_1 = \text{good} | S_0 = \text{good})P(S_0 = \text{good}) + P(S_1 = \text{good} | S_0 = \text{bad})P(S_0 = \text{bad}) .$$

Cool! We know those values! The conditional probabilities come from our sensor model and the others from our prior. So

$$\begin{aligned} P(S_1 = \text{good}) &= 0.7 \cdot 0.9 + 0.1 \cdot 0.1 \\ &= 0.64 \end{aligned}$$

Hmm. So, according to our transition model, these copy machines disintegrate pretty quickly! Without any new observations, we think that after the first time step, the machine only has probability 0.64 of being good.

Figure 2 shows a schematic version of this update rule, which is a good way to think about computing it either by hand or in a computer. To compute an entry in the distribution $P(S_1 = s_1)$, you take each element s_0 in $P(S_0 = s_0)$ and multiply it by the transition probability $P(S_1 = s_1 | S_0 = s_0)$, and then sum the results.

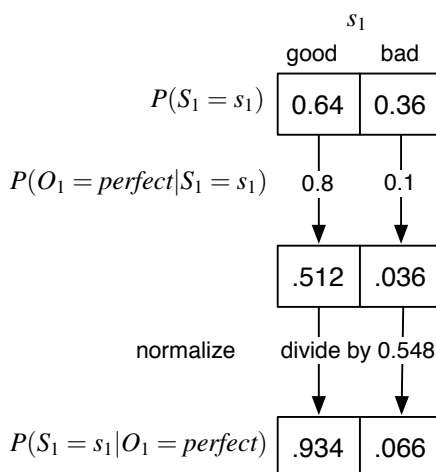


Figure 3: Schematic version of first observation update

Now it's time to take advantage of the information that it printed a perfect copy. We can use Bayes' rule to get:

$$P(S_1 = \text{good} | O_1 = \text{perfect}) = \frac{P(O_1 = \text{perfect} | S_1 = \text{good})P(S_1 = \text{good})}{P(O_1 = \text{perfect})}.$$

We know both of the terms in the numerator: one is the sensor model and the other is the value we just computed. But what is $P(O_1 = \text{perfect})$? We can derive it in a way similar to that shown above, reducing it to values that we've already calculated, or that we have in our sensor model.

$$\begin{aligned}
 P(O_1 = \text{perfect}) &= P(O_1 = \text{perfect}, S_1 = \text{good}) + P(O_1 = \text{perfect}, S_1 = \text{bad}) \\
 &= P(O_1 = \text{perfect} | S_1 = \text{good})P(S_1 = \text{good}) + P(O_1 = \text{perfect} | S_1 = \text{bad})P(S_1 = \text{bad}) \\
 &= 0.8 \cdot 0.64 + 0.1 \cdot 0.36 \\
 &= 0.548
 \end{aligned}$$

Now we can go back to our Bayes' rule expression:

$$\begin{aligned}
 P(S_1 = \text{good} | O_1 = \text{perfect}) &= \frac{P(O_1 = \text{perfect} | S_1 = \text{good})P(S_1 = \text{good})}{P(O_1 = \text{perfect})} \\
 &= \frac{0.8 \cdot 0.64}{0.548} \\
 &= 0.9343
 \end{aligned}$$

Figure 3 shows a schematic version of this update rule, which is a good way to think about computing it either by hand or in a computer. To compute an entry in the distribution $P(S_1 = s_1 | O_1 = \text{perfect})$, you take each element s_1 in $P(S_0 = s_1)$ and multiply it by the observation probability $P(O_1 = \text{perfect} | S_1 = s_1)$. Then, you need to normalize the distribution so that it sums to 1; so you divide each value by the sum of all the values.

Whew! After all that, we believe our copy machine more likely to be in a good state, than we did when we got it out of the box.

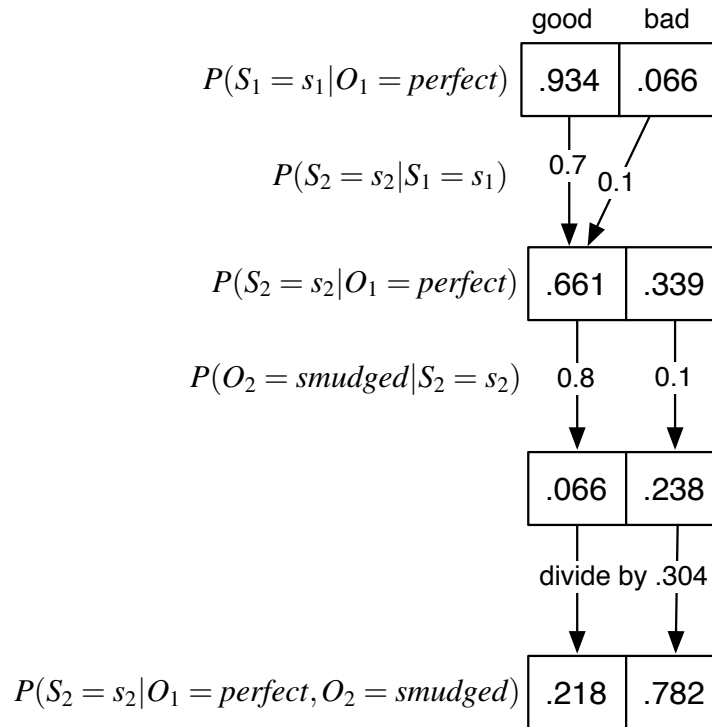


Figure 4: Schematic version of second transition and observation update

Our second copy Now, let's imagine we print another page, and it's smudged. We want to compute $P(S_2 = \text{good} | O_1 = \text{perfect}, O_2 = \text{smudged})$. Rather than writing out all the algebra again, we'll just do it in our schematic form, with both steps as shown in figure 4.

Ow. Now we're pretty sure our copy machine is no good. Planned obsolescence strikes again!

General state estimation

Just for completeness sake, we'll write out the state-update procedure for a general situation. Let \mathbf{b}_t be the *belief state* at time t , after incorporating actual observations $\mathbf{o}_1, \dots, \mathbf{o}_t$:

$$\mathbf{b}_t(s) = P(S_t = s | O_1 = \mathbf{o}_1 \dots O_t = \mathbf{o}_t) \quad .$$

Then we proceed in two steps:

Transition update :

$$\mathbf{b}'_{t+1}(s_{t+1}) = P(S_{t+1} = s_{t+1} | O_1 = \mathbf{o}_1 \dots O_t = \mathbf{o}_t) = \sum_{s_t} P(S_{t+1} = s_{t+1} | S_t = s_t) \mathbf{b}_t(s_t) \quad .$$

Observation update, given \mathbf{o}_{t+1} :

$$\mathbf{b}_{t+1}(s_{t+1}) = P(S_{t+1} = s_{t+1} | O_1 = \mathbf{o}_1 \dots O_{t+1} = \mathbf{o}_{t+1}) = \frac{P(O_{t+1} = \mathbf{o}_{t+1} | S_{t+1} = s_{t+1}) \mathbf{b}'_{t+1}(s_{t+1})}{\sum_{s_j} P(O_{t+1} = \mathbf{o}_{t+1} | S_{t+1} = s_j) \mathbf{b}'_{t+1}(s_j)} \quad .$$

A very important thing to see about these definitions is that they enable us to build what is known as a *recursive* state estimator. (Unfortunately, this is a different use of the term “recursive” than we’re used to from programming languages). It means that, after each action and observation, we can update our belief state, to get a new $\mathbf{b}_t(s)$. Then, we can forget the particular action and observation we had, and just use the $\mathbf{b}_t(s)$, \mathbf{a}_t , and \mathbf{o}_{t+1} to compute $\mathbf{b}_{t+1}(s)$.