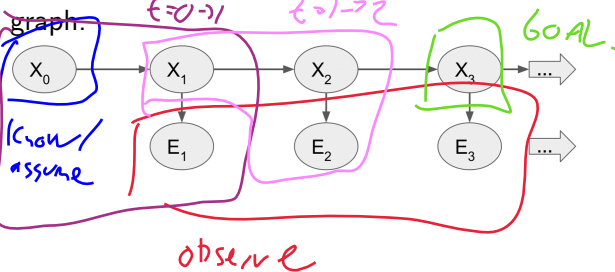


## Oct 26 HMM Wrapup and MDP

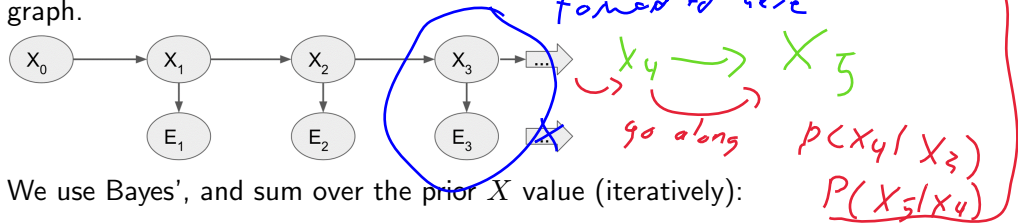
**Opening Example:** Sketch the FORWARD algorithm to find  $P(X_3|E_{1:3})$  on the given



What about  $P(X_5|E_{1:3})$ ?

## Oct 26 HMM Wrapup and MDP

**Opening Example:** Sketch the FORWARD algorithm to find  $P(X_3|E_{1:3})$  on the given graph.



We use Bayes', and sum over the prior  $X$  value (iteratively):

$$P(X|E) = \frac{P(E|X)P(X)}{P(E)}$$

Bayes thm.

$$\propto P(E|X) \sum_{\text{prior } X} P(X|\text{prior } X)P(\text{prior } X)$$

$P(X)$  based on prior steps

What about  $P(X_5|E_{1:3})$ ?

# Announcements and To-Dos

## Announcements:

1. Skip 1a for now but it's worth a bit of E.C. if you get  $A^*$  working. I'll add a few edges to hard code in an addendum.

Last time we learned:

1. Stationary distributions to Markov Models.

input (start node)  
(end node) ) disconnect

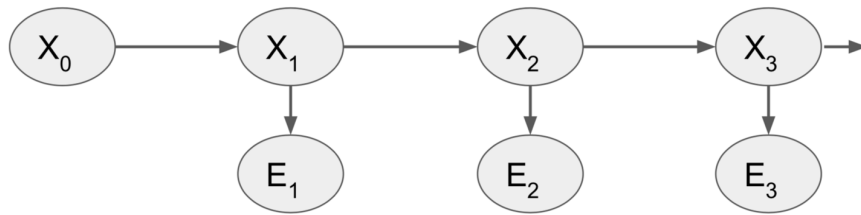
return shortest path

# Hidden Markov Models

## Example:

Suppose you are a graduate student in a basement office. You are writing your dissertation, so you don't get to leave very often.

You are curious if it is raining, and the only contact you have with the outside world is through your advisor. If it is raining, she brings her umbrella 90% of the time, and has it just in case on 20% of sunny days. You know that historically, 40% of rainy days were followed by another rainy day, and 30% of sunny days were followed by a rainy day.



## HMM: Filtering

**Filtering:** The goal is to predict  $X_{t+1}$  *given* all the evidence available  $E_{1:t+1}$ .

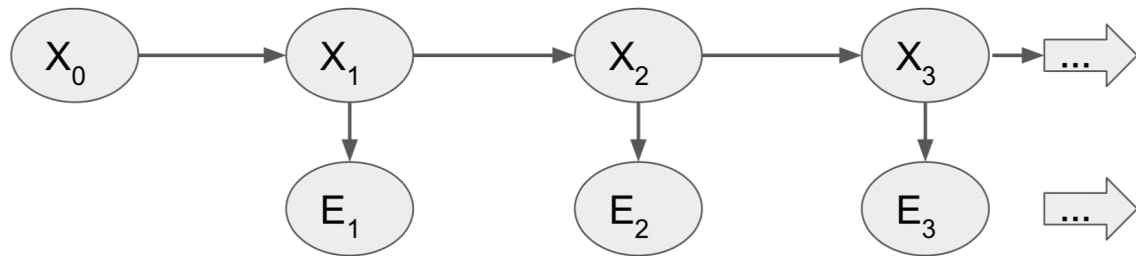
At  $t = 0$ :

$$P(X_1|E_1) = P(E_1|X_1) \sum_{X_0} P(X_1|X_0)P(X_0)$$

At  $t = 1$ :

$$P(X_2|E_{1:2}) = P(E_2|X_2) \sum_{X_1} P(X_2|X_1)P(X_1|E_1)$$

We continue FORWARD through the graph.

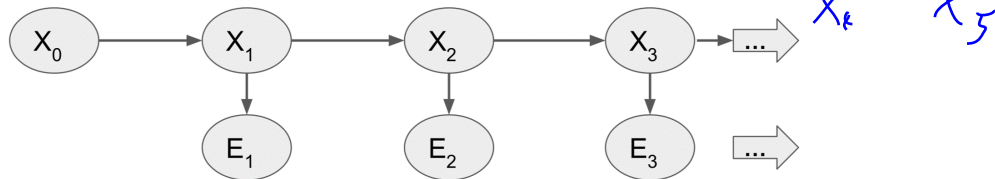


## HMM: Prediction

**Prediction:** The goal is to predict  $X_{t+k+1}$  *given* all the evidence available  $E_{1:t+1}$ .  
 The prediction of  $X$  two ( $k = 1$ ) time steps beyond where our evidence ended was:

$$P(X_{t+2}|E_{1:t}) = \underbrace{\sum_{X_{t+1}} P(X_{t+2}|X_{t+1})}_{\text{step of } t, \text{ Markov}} \underbrace{\sum_{X_t} P(X_{t+1}|X_t)}_{\text{Markov}} \underbrace{P(X_t|E_{1:t})}_{\text{from forward}}$$

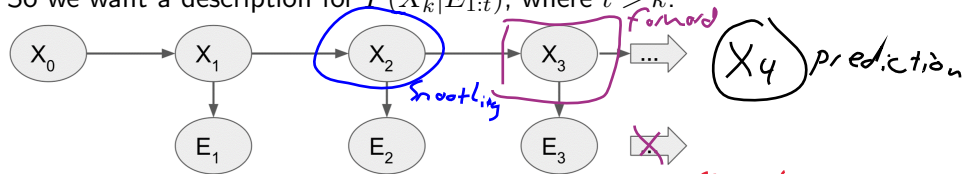
Making a  $k$ -step prediction just means doing FORWARD-steps up until we're out of evidence, and then following the Markov process to evolve  $X_{t+1}|X_t$  until we reach the desired future time.



## HMM: Smoothing

Our final task is **smoothing**, where we try to update probabilities of prior states  $X$  based on current evidence.

So we want a description for  $P(X_k|E_{1:t})$ , where  $t > k$ .



$$\begin{aligned}
 & \text{past} \quad \text{present} \quad \text{past} \quad \text{after } t. \\
 & \downarrow \quad \downarrow \\
 & P(X_k|E_{1:t}) = P(X_k|E_{1:k}, E_{k+1:t}) \\
 & P(X_k|E_{1:t}) = \alpha P(\underbrace{E_{k+1:t}}_{\text{past}} | \underbrace{X_k, E_{1:k}}_{\text{present}}) \underbrace{P(X_k|E_{1:k})}_{\text{past}} \quad \beta = \text{yes} \\
 & = \alpha P(E_{k+1:t} | X_k) P(X_k | E_{1:k})
 \end{aligned}$$

We can find the last term by the FORWARD algorithm for filtering.

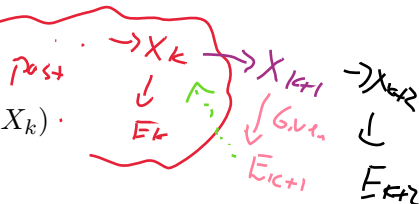
# HMM: Smoothing

This leaves the  $P(E_{k+1:t}|X_k)$  term, which we denote by  $b_{k+1:t}$ , which is the probability of future *measurements* given the current state of our system, which is just the combination of our transition and sensor models! Imagine taking *one* time step and asking about the new evidence: we need to describe  $X_{k+1}$ .

$$\begin{aligned}
 b_{k+1:t} &= P(E_{k+1:t}|X_k) \\
 &= \sum_{X_{k+1}} P(E_{k+1:t} | \underbrace{X_k, X_{k+1}}_{\text{indep}}) P(X_{k+1}, X_k) \\
 &= \sum_{X_{k+1}} \underbrace{P(E_{k+1:t} | X_{k+1})}_{\text{split up}} P(X_{k+1}, X_k) \\
 &= \sum_{X_{k+1}} \underbrace{P(E_{k+1}, E_{k+2:t} | X_{k+1})}_{\text{conditional}} P(X_{k+1}, X_k) \\
 &= \sum_{X_{k+1}} \underbrace{P(E_{k+1} | X_{k+1})}_{\text{sensor model}} \underbrace{P(E_{k+2:t} | X_{k+1})}_{\text{similar to LHS}} \underbrace{P(X_{k+1}, X_k)}_{\text{Markov model}}
 \end{aligned}$$

sum  
over the  
unobserved

missing  $X_{k+1}$





## HMM: Smoothing

The middle term is a *backwards* model, since we need the future value of  $P(E|X)$  rather than the past, like FORWARD did.

$$\begin{aligned}
 b_{k+1:t} &= P(E_{k+1:t}|X_k) \\
 &= \sum_{X_{k+1}} \underbrace{P(E_{k+1}|X_{k+1})}_{\text{sensor model}} \underbrace{P(E_{k+2:t}|X_{k+1})}_{\text{recursive}} \underbrace{P(X_{k+1}, X_k)}_{\text{Markov model}} \\
 &= \text{BACKWARD}(b_{k+2:t}, E_{k+1})
 \end{aligned}$$

All told, then, we have:

$$\begin{aligned}
 P(X_k|E_{1:t}) &= P(X_k|E_{1:k}, E_{k+1:t}) \\
 &\stackrel{\text{post } X}{=} \alpha \underbrace{P(E_{k+1:t}|X_k, E_{1:k})}_{\text{future } E} \underbrace{P(X_k|E_{1:k})}_{\text{forward to } k} \\
 &= \alpha \underbrace{P(E_{k+1:t}|X_k)}_{\text{Back to } k} P(X_k|E_{1:k}) \\
 &= \alpha \text{BACKWARD} \times \text{FORWARD} \\
 &= \alpha(b_{k+1:t}) \times (f_{1:k})
 \end{aligned}$$

# HMM: Smoothing

What does this actually look like?

$P(X_1|E_{1:3})$  = What's the probability it rained on Day 1 given 3 days of evidence (TFT)?

$$= \alpha f_{1:1} b_{2:3}$$

use of  $E_1 = T$  (green)  
 use of  $E_2 = F$  (blue)  
 use of  $E_3 = T$  (blue)

$f_{1:1}$  = What's the probability it rained on Day 1 given evidence through day 1?

$$= \underbrace{P(X_1|E_{1:1})}_{\text{time 0?}} = \alpha P(E_1|X_1) \sum_{X_0} P(X_1|X_0) P(X_0|E_{null})$$

$$= \begin{pmatrix} .708 \\ .292 \end{pmatrix}$$

# HMM: Smoothing

We have to run BACKWARDS for  $k = 1$  and  $k = 2$ . ( $t = 3$  for both!)

$b_{2:3}$  = What's the probability of the evidence on days 2 and 3, given  $X$  at day 2?

$$\begin{aligned}
 &= P(E_{2:3}|X_2) \\
 &= \sum_{X_2} P(E_2|X_2) \underbrace{P(E_{3:3}|X_2)}_{\text{evidence at time 3}} P(X_2|X_1) \\
 &= \sum_{X_2} P(E_2|X_2) b_{3:3} P(X_2|X_1)
 \end{aligned}$$

*Handwritten notes:* "evidence at time 2" with an arrow pointing to  $E_2$ ; "evidence at time 3" with a bracket around  $P(E_{3:3}|X_2)$ ; "tells us..." with an arrow pointing to the summation over  $X_2$ .

$b_{3:3}$  = What's the probability of the evidence on days 3-3, given  $X$  at day 2?

$$\begin{aligned}
 &= P(E_{3:3}|X_2) \\
 &= \sum_{X_3} P(E_3|X_3) P(E_{4:3}|X_3) P(X_3|X_2) \\
 &= \sum_{X_3} P(E_3|X_3) \underbrace{b_{4:3}}_{\text{last location of evidence}} P(X_3|X_2) \text{ but } b_{4:3} = 1 \text{ by independence!}
 \end{aligned}$$

*Handwritten notes:* "last location of evidence" with an arrow pointing to  $b_{4:3}$ .

# HMM: Smoothing

We have to run BACKWARDS for  $k = 1$  and  $k = 2$ . ( $t = 3$  for both!)

$$\begin{aligned}
 b_{2:3} &= \sum_{X_2} P(E_2|X_2)b_{3:3}P(X_2|X_1) \\
 b_{3:3} &= \sum_{X_3} \underbrace{P(E_3|X_3)}_{P(E_3=T|X_3)} \underbrace{P(X_3|X_2)}_{\text{green box}} \\
 &= \alpha \left[ \underbrace{\begin{pmatrix} 0.4 \\ 0.3 \end{pmatrix}}_{X_3=T} (0.9) + \underbrace{\begin{pmatrix} 0.6 \\ 0.7 \end{pmatrix}}_{X_3=F} (0.2) \right] = \begin{pmatrix} 0.48 \\ 0.41 \end{pmatrix}
 \end{aligned}$$

*Handwritten notes: Green boxes around  $b_{3:3}$  and the transition probability matrix. Red arrows point from  $P(E_3=T|X_3)$  to the  $0.9$  and  $0.2$  terms. Green arrows point from  $X_2=T$  and  $X_2=F$  to the corresponding columns of the matrix.*

## Sensor Model

$X_t$	$P(E_t X_t)$
T	.9
F	.2

## Transition Model

$X_t$	$P(X_{t+1} X_t)$
T	.4
F	.3

*Handwritten notes: Green circles around .4 and .3. Dashed green lines and numbers 6 and 7 to the right.*

## Initializations and Evidence

$$\begin{aligned}
 E_{1:3} &= [T, F, T] \\
 X_0 &= [.5, .5]
 \end{aligned}$$

*Handwritten notes: Blue circle around the last 'T' in the evidence sequence.*

# HMM: Smoothing

We have to run BACKWARDS for  $k = 1$  and  $k = 2$ . ( $t = 3$  for both!)

$$b_{2:3} = \sum_{X_2} P(E_2|X_2)b_{3:3}P(X_2|X_1)$$

*b<sub>3:3</sub> top*                      *b<sub>3:3</sub> bottom*

$$= \alpha \left[ \underbrace{\begin{pmatrix} 0.4 \\ 0.3 \end{pmatrix} \begin{pmatrix} 0.48 \\ 0.1 \end{pmatrix}}_{X_2=T} + \underbrace{\begin{pmatrix} 0.6 \\ 0.7 \end{pmatrix} \begin{pmatrix} 0.41 \\ 0.8 \end{pmatrix}}_{X_2=F} \right] = \begin{pmatrix} 0.68 \\ 0.32 \end{pmatrix}$$

*P(E<sub>2</sub>=T)*

**Sensor Model**

$X_t$	$P(E_t X_t)$
T	.9
F	.2

*E<sub>2</sub>=F*  
*.1*  
*.8*

**Transition Model**

$X_t$	$P(X_{t+1} X_t)$
T	.4
F	.3

**Initializations and Evidence**

$E_{1:3} = [T, F, T]$ ,  
 $X_0 = [.5, .5]$

## HMM: Smoothing

## Sensor Model

$X_t$	$P(E_t X_t)$
T	.9
F	.2

## Transition Model

$X_t$	$P(X_{t+1} X_t)$
T	.4
F	.3

## Initializations and Evidence

$$E_{1:3} = [T, F, T],$$

$$X_0 = [.5, .5]$$

We had two calculations:

$$P(X_1|E_1) = \begin{pmatrix} 0.708 \\ 0.292 \end{pmatrix} \quad P(X_1|E_{1:3}) = \begin{pmatrix} 0.682 \\ 0.318 \end{pmatrix}$$

Sanity check? Why is the  $P(X_1 = T)$  smaller with more evidence?

Day 2:  
Evidence =  
no umbrella

... maybe sun  
on day 2  
... maybe sun  
on day 1?

# HMM: Smoothing

## Sensor Model

$X_t$	$P(E_t X_t)$
T	.9
F	.2

## Transition Model

$X_t$	$P(X_{t+1} X_t)$
T	.4
F	.3

## Initializations and Evidence

$$E_{1:3} = [T, F, T],$$

$$X_0 = [.5, .5]$$

We had two calculations:

$$P(X_1|E_1) = \begin{pmatrix} 0.708 \\ 0.292 \end{pmatrix} \quad P(X_1|E_{1:3}) = \begin{pmatrix} 0.682 \\ 0.318 \end{pmatrix}$$

Sanity check? Why is the  $P(X_1 = T)$  smaller with more evidence?

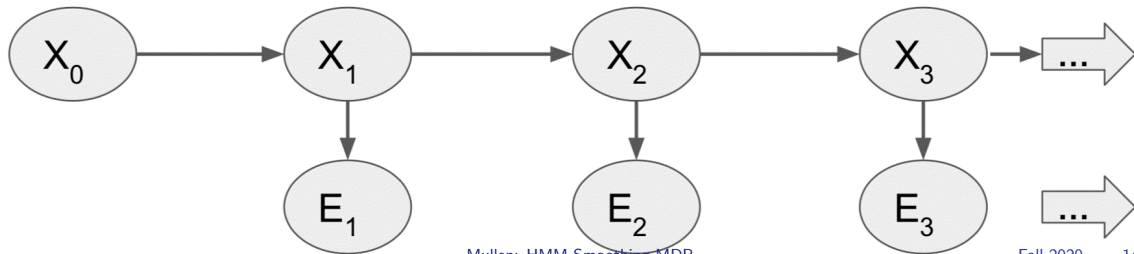
**Solution:** we saw evidence on rain on day 2, which in turn means it was likely to have rained on day 2... and a rainy day more likely preceded by another rainy one!

## HMM: All at Once

So for any given observation  $X_k$ , we tend to have to run both a forward algorithm to ask what the evidence *up to time  $k$*  did, then a backwards algorithm to ask what the evidence afterwards did. To solve the whole chain, we do *both*. Given evidence up to time  $t$ , we:

- ▶ Run the FORWARD algorithm to filter it.
- ▶ then run the BACKWARD algorithm to smooth it

We use  $f_{1:k}$  in the backwards algorithm, so we'll save them: the main tenet of dynamic programming is to not solve the same problem twice!





## HMM: Wrapup

There's a final question that often is asked: what's the *most likely* sequence of  $X$  values that gave rise to our evidence.

- ▶ Lazy way: compute the  $P(X|E)$  values and pick the most likely one for each time individually.
- ▶ Rigorous way: compute a maximization over all the nodes of  $P(X_0, X_1, \dots, X_t, E_0, E_1, \dots, E_t)$

It turns out the rigorous problem can heavily exploit our independence assumptions, as usual! The joint density of the HMM will factor into

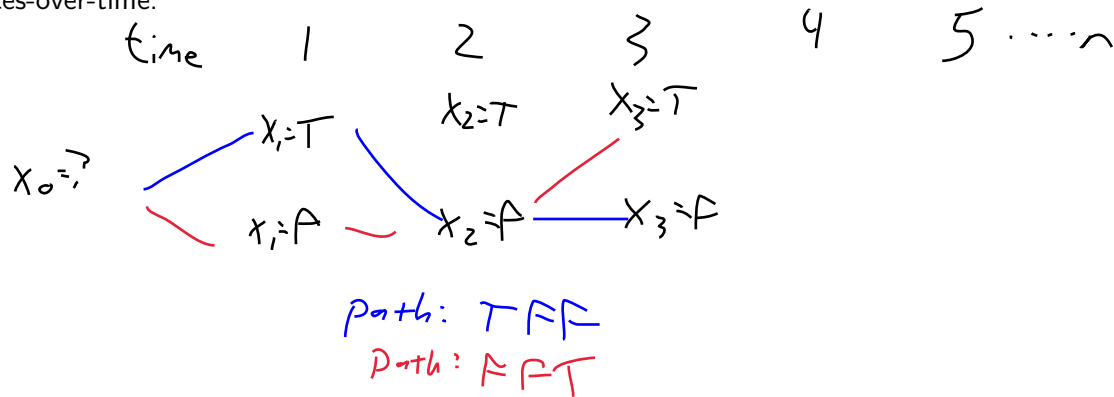
$$\begin{aligned} & \prod_{all\ nodes} P(Z_i | \text{parents}(Z_i)) \\ &= P(X_0)P(X_1|X_0)P(E_1|X_1)P(X_2|X_1)P(E_2|X_2) \dots \end{aligned}$$

## HMM: Paths

AND

$$P(X_0)P(X_1|X_0)P(E_1|X_1)P(X_2|X_1)P(E_2|X_2)\dots$$

This is the probability of a specific sequence or *path* through the graph of  $X$  states-over-time.



## HMM: Viterbi Overview

The recursive algorithm for this is called the *Viterbi* algorithm. But let's think about the calculation by hand, briefly. Suppose we have an initial probability of  $P(rain) = 1$ . Then we observe over 2 days: *Umbrella, None*.

$$E_1 = T \quad E_2 = F$$

$$X_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

**Sensor Model**

$X_t$	$P(E_t X_t)$
T	.9
F	.2

**Transition Model**

$X_t$	$P(X_{t+1} X_t)$
T	.4
F	.3

**Initializations**

$$E_{1:2} = [T, F], X_0 = [1, 0]$$

Our task now is to compute and compare:

1.  $P(X_1 = rain|E)$  vs.  $P(X_1 = sun|E)$

forward to  $X_1$

2.  $P(X_2 = rain|E)$  vs.  $P(X_2 = sun|E)$

" to  $X_2$

3. The full chains where one of (1) is true and one of (2) is true. We have to compare all of  $RR, RS, SR, SS$  and choose the *most likely* sequence.

4 chains of length 2.

## HMM: Viterbi Overview

Sensor Model	
$X_t$	$P(E_t X_t)$
T	.9
F	.2

Transition Model	
$X_t$	$P(X_{t+1} X_t)$
T	.4
F	.3

## Initializations

$$E_{1:2} = [T, \cancel{F}], X_0 = [1, 0]$$

We can compute two probabilities at time  $t = 1$ , just as in forward-stepping.

1. Joint probability that evidence at  $t = 1$  was  $U$  and the true value was *rain*. This is

$$\underline{P(X_1 = T, E_1 = T)} = \underline{P(E_1 = T|X_1 = T)} \underline{P(X_1 = T)}$$

$$P(X_1 = T, E_1 = T) = P(E_1 = T|X_1 = T) \sum_{X_0} P(X_1 = T|X_0)$$

*split / sum over*

2. Joint probability that evidence at  $t = 1$  was  $U$  and the true value was *sun*. This is

$$P(X_1 = F, E_1 = T) = P(E_1 = T|X_1 = F)P(X_1 = F)$$

$$P(X_1 = F, E_1 = T) = P(E_1 = T|X_1 = F) \sum_{X_0} P(X_1 = F|X_0)$$

## HMM: Viterbi Overview

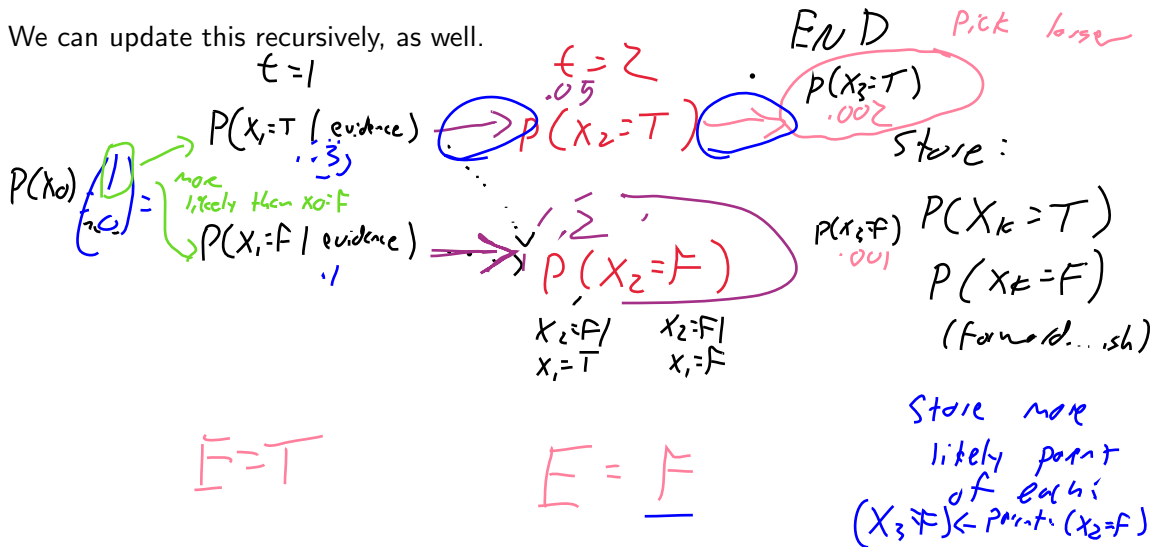
We then *classify* what was the most likely outcome at the first time step. This is also a way to ignore the denominator ( $\alpha$ ) from Bayes that occurred in forward-stepping.

The key to the Viterbi is to then use the joint probabilities through time  $t = 1$  to compute probabilities up to time  $t = 2$ . There are two ways to get a “rain” at time  $t = 2$ : the ones that came from  $X_1 = \text{rain}$  and the ones that came from  $X_1 = \text{sun}$ , and we previously computed each of those.

Writing down the full probability for  $P(X_2 = \text{rain})$  is messy, but instead we just break down using the cases on  $X_1$  we already have.

## HMM: Viterbi Sketch

We can update this recursively, as well.



## HMM: Most Likely Sequence

$$P(X, E) = P(X_0)P(X_1|X_0)P(E_1|X_1)P(X_2|X_1)P(E_2|X_2) \dots$$

So we want to maximize this thing... but maximizing products is harder than sums, so we hit with a log for numerical stability, which keeps the max in the same place and changes products to sums.

$$\log P(X, E) = \log (P(X_0)P(X_1|X_0)P(E_1|X_1)) + \sum \log (\underbrace{P(X_k|X_{k-1})}_{\text{transition}} \underbrace{P(E_k|X_k)}_{\text{emission}})$$

The recursive algorithm for this is called the *Viterbi* algorithm and is computed in linear time. It's just the calculations above where, for each time step and for each state (T/F, umbrella or not), we store the summed up probabilities of "all the ways we could get to this state."

To choose the best path at the end, we just work backwards. We choose the best *state* at the end, and then go backwards along the graph asking what that states' most likely predecessor was.

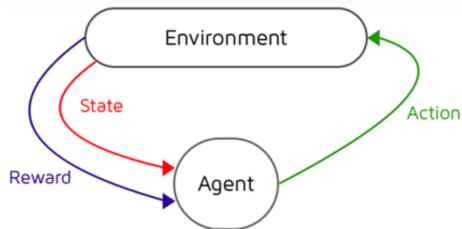
## Markov Decision Process

A *Markov Decision Process* (MDP) is a sequential decision problem that is the combination of a Markov Model and a decision-making agent. It asks the question: how do we maximize utility if there are uncertainties associated with the successor states of each action? To do this, we require:

1. A fully observable, stochastic environment
2. A Markov transition model that gives probabilities of states *given* decisions
3. An additive reward structure

They are often used for

1. Inventory management
2. Routing/logistics
3. Games
4. Planning under uncertainty





## MDPs

Consider an agent-based game. We win if we reach the treasure. We lose if we run into the internet troll or the dragon.

**Goal:** describe the appropriate set of moves from our current location to the treasure.



## MDPs

Consider an agent-based game. We win if we reach the treasure. We lose if we run into the internet troll or the dragon.

**Goal:** describe the appropriate set of moves from our current location to the treasure.

Suppose the Dragon and Troll can't move.  
What do we do?



## MDPs

Consider an agent-based game. We win if we reach the treasure. We lose if we run into the internet troll or the dragon.

**Goal:** describe the appropriate set of moves from our current location to the treasure.

Suppose the Dragon and Troll can't move.  
What do we do?

Twist: suppose, as in our trip to Taco Bell, we sometimes get disoriented, and move in a *different* direction than the one we choose!



## MDP Uncertainty and Choice

The MDP is meant to describe a real world process where actions are not perfectly reliable. Suppose we describe a transition model:

1. *Given* our intended action, the probability we move where we intend is .8.
2. *Given* our intended action, the probability we move to either side ( $90^\circ$ ) is .1 each.

**Definition:** A *policy* is what we would tell our agent to do in any given possible state  $s$ .

Denote  $\pi(s)$  by the policy chosen at state  $s$ .



## MDP Rewards

The MDP is meant to describe a real world process where actions are not perfectly reliable. Suppose we describe a transition model:

1. *Given* our intended action, the probability we move where we intend is .8.
2. *Given* our intended action, the probability we move to either side ( $90^\circ$ ) is .1 each.

**Definition:** A *policy* is what we would tell our agent to do in any given possible state  $s$ .

Denote  $\pi(s)$  by the policy chosen at state  $s$ .



## MDP Rewards

To choose a policy we need a notion of what makes a move good or bad. Suppose that:

1. Moving to the dragon or troll achieves a reward of -1, and ends the game.
2. Moving to the treasure achieves a reward of 1, and ends the game.
3. We can define a reward  $R(s)$  (or maybe  $R(s \rightarrow s')$ ) associated with moving to any state.

We may even encode a reward for the non-movement  $s \rightarrow s$ . For example,  $R(s \rightarrow s) > 0$ , an agent will rarely move, whereas a reward of  $R(s \rightarrow s) = -2$  will create a frenetic, always-moving agent.



## MDP Utility

Since rewards may not exist on all actions, we need to conceptualize an *expected* rewards or a *long-run* rewards. These like in the *utility* associated with each state.

**Utility** is our long-run gain.

1. It depends on the entire *sequence* of states visited,  $[s_0, s_1, \dots s_{50}, s_{51}, \dots]$
2. Informal definition: utility is the *sum* of rewards achieved over a set of states/movements:

$$U[s_0, s_1, \dots s_{50}, s_{51}, \dots] = R(s_0) + R(s_1) + \dots$$

Classically, two terms are added to clarify and add tuning to the concept of utility.



## MDP Utility

**Definition:** The *Time Horizon* of an MDP can be either:

1. *Finite Horizon*, where after a fixed time  $N$  no actions matter. Here we consider the rewards or utility  $U[s_0, s_1, \dots s_N, s_{N+1}, \dots] = U[s_0, s_1, \dots s_N]$ . The length of the horizon may impact your decisions.

**Example:** It's the first/last lap of your game of Mario Kart. Should you save your star piece for when someone tries to shoot you?

2. *Infinite Horizon*, where there is never a reason to behave differently in the same state at different times.

**Example:** When you get the treasure doesn't matter, only whether you get there.

**Definition:** The *discount factor* of an MDP is a multiplicative punishment  $\gamma$  for taking longer to reach rewards. It's common in finance as it represents an increased value of immediate rewards over future rewards.

$$U[s_0, s_1, \dots s_{50}, s_{51}, \dots] = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \gamma^3 R(s_3) + \dots$$



## MDP Goal

So we have:

1. A Markov chain that gives successors *given* actions
2. Rewards associated with each state
3. A utility that may track rewards of a *sequence* of states
4. A possible discount on *when* we get rewards
5. A preference for how many total moves matter

**Goal:** The output of an MDP is an *optimal policy* that specifies where to move from a given starting state  $s$ . We maximize the expected utility under policy  $\pi$ :

$$E[U^\pi(s)] = E\left[\sum_{t=0}^{\infty} \gamma^t R(S_t)\right]$$

## MDP Rewards

Suppose we start at the state  $(3, 1)$ . What is the *expected utility* of the policy of “move to the right”?

At time  $t = 0$ , we have utility of  $\gamma^0 R((3, 1))$ .

1. 80% chance we actually go right and achieve a reward of  $+1$ , for utility  $\gamma^1$ .
2. 10% chance we actually go up and achieve a reward of  $-1$ , for utility  $-\gamma^1$ .
3. 10% chance we actually go right and achieve a reward of... whatever the discounted *utility* of tile  $(2, 1)$  is.



## MDP Rewards

Suppose we start at the state  $(3, 1)$ . What is the *expected utility* of the policy of “move to the right”?

At time  $t = 0$ , we have utility of  $\gamma^0 R((3, 1))$ .

1. 80% chance we actually go right and achieve a reward of  $+1$ , for utility  $\gamma^1$ .
2. 10% chance we actually go up and achieve a reward of  $-1$ , for utility  $-\gamma^1$ .
3. 10% chance we actually go right and achieve a reward of... whatever the discounted *utility* of tile  $(2, 1)$  is.



## MDP Rewards

Suppose we start at the state  $(3, 1)$ . What is the *expected utility* of the policy of “move to the right”?

At time  $t = 0$ , we have utility of  $\gamma^0 R((3, 1))$ .

1. 80% chance we actually go right and achieve a reward of  $+1$ , for utility  $\gamma^1$ .
2. 10% chance we actually go up and achieve a reward of  $-1$ , for utility  $-\gamma^1$ .
3. 10% chance we actually go right and achieve a reward of... whatever the discounted *utility* of tile  $(2, 1)$  is.

So at  $t = 1$ ,  $U^\pi((3, 1)) =$

$$\gamma [0.1R(3, 2) + 0.8R(4, 1) + 0.1R(3, 1)]$$



## MDP Rewards: What to consider

Our utility gained after one attempt to move right was

$$\begin{aligned} U^\pi((3,1)) &= \gamma [.1R(3,2) + .8R(4,1) + .1R(3,1)] \\ &= \gamma [.1(-1) + .8(1) + .1R(3,1)] \end{aligned}$$

If we're allowed to take more moves, the  $R(3,1)$  term should get it's own policy: where should we move from  $(3,1)$  if it's currently  $t = 1$ ?

But the value of  $U^\pi((3,1))$  at  $t = 1$  isn't necessarily the same as the left-hand side of  $U^\pi((3,1))$  at  $t = 0$ ! This now depends on our *horizon*. If we only had one more left, it might be utility of zero!



## MDP Algorithm:

Our goal with an MDP is to maximize the expected discounted utility after playing the game to its horizon. So we compute the utility associated with any given policy  $\pi$  and choose the best one, the *optimal policy*  $\pi^*(s)$ :

$$\pi^*(s) = \arg \max_{\pi} U^{\pi}(s)$$

**Definition:** The *true utility* of a state is its utility when associated with the optimal policy  $\pi^*$ . We denote it  $U^{\pi^*}(s)$  or just  $U(s)$ .

**Result:** The *true utility* of a state is the expected utility gained by choosing the best successor state. This is the sum of the discounted utilities of all the possible successors of state  $s$  under optimal decision  $a$ .

# Moving Forward

- ▶ Coming up:
  1. Computing stuff on Markov Decision Processes!
  2. Markov NB on Friday.