

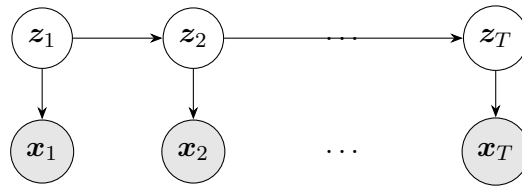
# CS 1810 Spring 2025 Section 9 Notes:

## HMMs and MDPs

### 1 Hidden Markov Models

A Hidden Markov Model (HMM) is useful for inferring a sequence of unknown or hidden states from a corresponding sequence of observed evidence.

#### 1.1 Graphical Model and Properties



We have a process that can be in one of  $K$  possible *states*,  $C_1, \dots, C_K$ , at each timestep. The state at a given timestep  $t$  is captured by the latent variable  $z_t$ . States cannot be observed, but they generate *observations*  $x_t$  which we can see. There are  $M$  possible values,  $O_1, \dots, O_M$ , that each observation can take on. The state the model is in at time  $t+1$  only depends on the state it was in at time  $t$ . A full HMM chain consists of  $T$  observations. Note that we then define a dataset in the context of HMM's as a collection of  $N$  of these chains, so that the dataset is written as  $\{\mathbf{x}^{(n)}\}_{n=1}^N$  where  $\mathbf{x}^{(n)} = (x_1^{(n)}, \dots, x_T^{(n)})$ .

The key properties of HMMs are listed below.

- (*Markov property*) The next hidden state depends only on the current hidden state and nothing else.

$$p(z_{t+1} \mid z_1, \dots, z_t, x_1, \dots, x_t) = p(z_{t+1} \mid z_t),$$

for  $t = 1, 2, \dots, T-1$ .

- The current observation depends only on the current hidden state.

$$p(x_t \mid z_1, \dots, z_t, x_1, \dots, x_{t-1}) = p(x_t \mid z_t),$$

for  $t = 1, 2, \dots, T$ .

Observe that we can read these properties off of the graphical model!

#### 1.2 Exercise: When to Use HMMs (Source: CMU)

For each of the following scenarios, is it appropriate to use a Hidden Markov Model? Why or why not? What would the observed data be in each case, and what would the hidden states capture?

1. Stock market price data
2. Recommendations on a database of movie reviews
3. Daily precipitation data in Boston
4. Optical character recognition for identifying words

### 1.3 Parameterization

We can use the HMM assumptions from above to decompose the joint distribution of the hidden states and observations into a more useful form. First, note that

$$p(\mathbf{z}_1, \dots, \mathbf{z}_T, \mathbf{x}_1, \dots, \mathbf{x}_T) = p(\mathbf{z}_1, \dots, \mathbf{z}_T) p(\mathbf{x}_1, \dots, \mathbf{x}_T \mid \mathbf{z}_1, \dots, \mathbf{z}_T)$$

Now for the left term on the RHS, we have

$$\begin{aligned} p(\mathbf{z}_1, \dots, \mathbf{z}_T) &= p(\mathbf{z}_1, \dots, \mathbf{z}_{T-1}) p(\mathbf{z}_T \mid \mathbf{z}_1, \dots, \mathbf{z}_{T-1}) \\ &= p(\mathbf{z}_1, \dots, \mathbf{z}_{T-1}) p(\mathbf{z}_T \mid \mathbf{z}_{T-1}) \\ &\quad \vdots \\ &= p(\mathbf{z}_1) \prod_{t=1}^{T-1} p(\mathbf{z}_{t+1} \mid \mathbf{z}_t) \end{aligned}$$

As for the right term, we have

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_T \mid \mathbf{z}_1, \dots, \mathbf{z}_T) &= p(\mathbf{x}_1, \dots, \mathbf{x}_{T-1} \mid \mathbf{z}_1, \dots, \mathbf{z}_T) p(\mathbf{x}_T \mid \mathbf{z}_1, \dots, \mathbf{z}_T, \mathbf{x}_1, \dots, \mathbf{x}_{T-1}) \\ &= p(\mathbf{x}_1, \dots, \mathbf{x}_{T-1} \mid \mathbf{z}_1, \dots, \mathbf{z}_{T-1}) p(\mathbf{x}_T \mid \mathbf{z}_T) \\ &\quad \vdots \\ &= \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{z}_t) \end{aligned}$$

In the second equality we use the fact that

$$p(\mathbf{x}_1, \dots, \mathbf{x}_{T-1} \mid \mathbf{z}_1, \dots, \mathbf{z}_T) = p(\mathbf{x}_1, \dots, \mathbf{x}_{T-1} \mid \mathbf{z}_1, \dots, \mathbf{z}_{T-1}),$$

which is true because  $\mathbf{z}_{T-1}$  blocks all paths between  $\mathbf{z}_T$  and all the observations  $\mathbf{x}_1, \dots, \mathbf{x}_{T-1}$ . Putting these equalities together, it follows that

$$p(\mathbf{z}_1, \dots, \mathbf{z}_T, \mathbf{x}_1, \dots, \mathbf{x}_T) = p(\mathbf{z}_1) \prod_{t=1}^{T-1} p(\mathbf{z}_{t+1} \mid \mathbf{z}_t) \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{z}_t)$$

From this decomposition, we see that we need three parameters to fully specify our HMM:

1.  $\boldsymbol{\theta} \in \mathbb{R}^K$ : defines the prior distribution over initial hidden states  $\mathbf{z}_1$ . This corresponds to the term  $p(\mathbf{z}_1)$ .
2.  $\mathbf{T} \in \mathbb{R}^{K \times K}$ : transition matrix where  $T_{ij}$  is the probability of transitioning from latent state  $C_i$  to latent state  $C_j$ . This corresponds to terms of the form  $p(\mathbf{z}_{t+1} | \mathbf{z}_t)$ .
3.  $\boldsymbol{\pi} \in \mathbb{R}^{K \times M}$ : conditional probability matrix where  $\pi_{kl}$  is the probability of observing  $O_l$  given the latent state  $C_k$ . This corresponds to terms of the form  $p(\mathbf{x}_t | \mathbf{z}_t)$ .

To learn these parameters, we use expectation maximization.

## 1.4 Forward-Backward Algorithm

The HMM model is characterized by the joint distribution  $p(\mathbf{z}_1, \dots, \mathbf{z}_T, \mathbf{x}_1, \dots, \mathbf{x}_T)$ , which means that many of our training and inference tasks require marginalization to obtain conditionals. Thus, naive algorithms can be expensive (they require lots of nested summations over states), and we use EM instead. To compute the probabilities that we use in the E-step, we rely on two key quantities:  $\alpha_t(\mathbf{z}_t)$  and  $\beta_t(\mathbf{z}_t)$ . Below are the definitions and intuitions behind these two quantities:

- $\alpha_t(\mathbf{z}_t)$  represents the joint probability of observations  $\mathbf{x}_1, \dots, \mathbf{x}_t$  and state  $\mathbf{z}_t$ :

$$\begin{aligned}
\alpha_t(\mathbf{z}_t) &= p(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{z}_t) \\
&= p(\mathbf{x}_t | \mathbf{z}_t, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) p(\mathbf{z}_t, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) \\
&= p(\mathbf{x}_t | \mathbf{z}_t) \sum_{\mathbf{z}_{t-1}} p(\mathbf{z}_t, \mathbf{z}_{t-1}, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) \\
&= p(\mathbf{x}_t | \mathbf{z}_t) \sum_{\mathbf{z}_{t-1}} p(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) p(\mathbf{z}_{t-1}, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) \\
&= p(\mathbf{x}_t | \mathbf{z}_t) \sum_{\mathbf{z}_{t-1}} p(\mathbf{z}_t | \mathbf{z}_{t-1}) \alpha_{t-1}(\mathbf{z}_{t-1})
\end{aligned}$$

Thus, we have that

$$\alpha_t(\mathbf{z}_t) = \begin{cases} p(\mathbf{x}_1 | \mathbf{z}_1) p(\mathbf{z}_1) & \text{if } t = 1 \\ p(\mathbf{x}_t | \mathbf{z}_t) \sum_{\mathbf{z}_{t-1}} p(\mathbf{z}_t | \mathbf{z}_{t-1}) \alpha_{t-1}(\mathbf{z}_{t-1}) & \text{if } 1 < t \leq T \end{cases}$$

Intuitively, this allows us to quickly answer the question: “how likely are we to currently be in state  $\mathbf{z}_t$ , if we observed a specific list of values?” As we see above, it turns out that  $\alpha_t$  can be defined in terms of  $\alpha_{t-1}$ . That is, we move **forwards** through the sequence to calculate the  $\alpha$ ’s.

- $\beta_t(\mathbf{z}_t)$  represents the joint probability of observations  $\mathbf{x}_{t+1}, \dots, \mathbf{x}_T$  conditioned on state  $\mathbf{z}_t$ :

$$\begin{aligned}
\beta_t(\mathbf{z}_t) &= p(\mathbf{x}_{t+1}, \dots, \mathbf{x}_T | \mathbf{z}_t) \\
&= \sum_{\mathbf{z}_{t+1}} p(\mathbf{x}_{t+1}, \dots, \mathbf{x}_T, \mathbf{z}_{t+1} | \mathbf{z}_t) \\
&= \sum_{\mathbf{z}_{t+1}} p(\mathbf{x}_{t+1} | \mathbf{z}_t, \mathbf{z}_{t+1}, \mathbf{x}_{t+2}, \dots, \mathbf{x}_T) p(\mathbf{z}_{t+1}, \mathbf{x}_{t+2}, \dots, \mathbf{x}_T | \mathbf{z}_t) \\
&= \sum_{\mathbf{z}_{t+1}} p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) p(\mathbf{x}_{t+2}, \dots, \mathbf{x}_T | \mathbf{z}_{t+1}, \mathbf{z}_t) p(\mathbf{z}_{t+1} | \mathbf{z}_t) \\
&= \sum_{\mathbf{z}_{t+1}} p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) \beta_{t+1}(\mathbf{z}_{t+1})
\end{aligned}$$

Thus, we have that

$$\beta_t(\mathbf{z}_t) = \begin{cases} 1 & \text{if } t = T \\ \sum_{\mathbf{z}_{t+1}} p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) \beta_{t+1}(\mathbf{z}_{t+1}) & \text{if } 1 \leq t < T \end{cases}$$

We can think about this intuitively by asking “what are the chances of the next observations if we are currently in state  $\mathbf{z}_t$ ?” It turns out that  $\beta_t$  can be defined in terms of  $\beta_{t+1}$ . That is, we move **backwards** through the sequence to calculate the  $\beta$ 's.

Note that the probabilities we use for calculating  $\alpha$  and  $\beta$  are given by the parameters that we fix in the E-Step.

## 1.5 EM for HMMs

Given data points  $\{\mathbf{x}^{(n)}\}_{n=1}^N$  defined by sequences  $(x_1^{(n)}, \dots, x_T^{(n)})$  of length  $T$  represented as row vectors, we want to infer the parameters  $\{\mathbf{T}, \boldsymbol{\theta}, \boldsymbol{\pi}\}$ . Had we been given the true states, we could easily compute joint probability  $p(\mathbf{x}^{(n)}, \mathbf{z}^{(n)})$  and write the complete-data log likelihood, and maximize with respect to the parameters. Instead, we need to estimate state distributions and parameters iteratively.

### 1.5.1 Inference Patterns with $\alpha, \beta$

The following patterns are useful for inference with a trained HMM as well as during the E-Step:

- **Filtration:**  $p(\mathbf{z}_t | \mathbf{x}_1, \dots, \mathbf{x}_t) \propto p(\mathbf{z}_t, \mathbf{x}_1, \dots, \mathbf{x}_t) = \alpha_t(\mathbf{z}_t)$
- **Smoothing:**  $p(\mathbf{z}_t | \mathbf{x}_1, \dots, \mathbf{x}_T) \propto p(\mathbf{x}_1, \dots, \mathbf{x}_T, \mathbf{z}_t) = \alpha_t(\mathbf{z}_t) \beta_t(\mathbf{z}_t)$
- **Prediction:**  $p(\mathbf{x}_{T+1} | \mathbf{x}_1, \dots, \mathbf{x}_T) \propto \sum_{\mathbf{z}_T, \mathbf{z}_{T+1}} \alpha_T(\mathbf{z}_T) p(\mathbf{z}_{T+1} | \mathbf{z}_T) p(\mathbf{x}_{T+1} | \mathbf{z}_{T+1})$
- **Transition:**  $p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{x}_1, \dots, \mathbf{x}_T) \propto \alpha_t(\mathbf{z}_t) p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) \beta_{t+1}(\mathbf{z}_{t+1})$
- Joint of observations:  $p(\mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{\mathbf{z}_t} \alpha_t(\mathbf{z}_t) \beta_t(\mathbf{z}_t)$  (for any  $t$ )

The derivations of the above equations can be found in the textbook. As we can see above, the  $\alpha$ 's and  $\beta$ 's allow us to concisely capture the quantities we would like to model, which allows us to quickly compute important quantities.

### 1.5.2 E-Step

The goal of the expectation step is to compute the expected values of the hidden states given a fixed set of parameters  $\mathbf{w} = \{\mathbf{T}, \boldsymbol{\theta}, \boldsymbol{\pi}\}$ . That is, we estimate the state distribution for  $p(\mathbf{z}_1^{(n)}, \dots, \mathbf{z}_T^{(n)} | \mathbf{x}^{(n)})$ , which we will call  $\mathbf{q}^{(n)}$ . Note that  $\mathbf{x}^{(n)}$  contains all  $T$  timesteps, so  $\mathbf{q}^{(n)}$  is a matrix with  $T$  rows and  $K$  columns, where the rows correspond to each  $\mathbf{z}_j$ , and the columns correspond to the possible classes of the  $\mathbf{z}_j$ .

The idea is to find successive approximations of this quantity based on the data we have available. We let  $z_{t,k}^{(n)}$  be the indicator that  $\mathbf{z}_t = C_k$ . Then we define the  $t, k$  element of  $\mathbf{q}$  to be

$$q_{t,k}^{(n)} = E[z_{t,k}^{(n)} | \mathbf{x}^{(n)}] = P(\mathbf{z}_t^{(n)} = C_k | \mathbf{x}^{(n)}).$$

That is, this is the probability that the state at time  $t$  is in class  $k$  given all the observed emissions. Notice how this is exactly the *smoothing* quantity we had in the previous subsection, which is the motivation for defining  $\alpha_t$  and  $\beta_t$ .

We would also need to consider the expectation of the *joint* of two consecutive states. Mathematically, this is written as  $\mathbf{Q}_{t,t+1}^{(n)} = E[\mathbf{z}_t^{(n)}, \mathbf{z}_{t+1}^{(n)} | \mathbf{x}^{(n)}]$ . Note that to encapsulate all possible values of the states, this would mean that  $\mathbf{Q}_{t,t+1}^{(n)}$  is a matrix. We then define the  $k, l$  element to be

$$Q_{t,t+1,k,l}^{(n)} = E[z_{t,k}^{(n)}, z_{t+1,l}^{(n)} | \mathbf{x}^{(n)}] = P(\mathbf{z}_t^{(n)} = C_k, \mathbf{z}_{t+1}^{(n)} = C_l | \mathbf{x}^{(n)}).$$

Notice how this is exactly the transition equation in the previous subsection!

### 1.5.3 M-Step

Now we need to update our parameters to maximize the expected complete-data log likelihood  $\mathbb{E}_{\mathbf{z}}[\ln p(\mathbf{x}, \mathbf{z}; \mathbf{w})]$ . It becomes a very nasty expression, so we will not include it here. It can be found in the textbook for reference. Applying the appropriate Lagrange multipliers and maximizing with respect to each of the parameters of interest, we recover the following update equations:

$$\theta_k = \frac{\sum_{n=1}^N q_{1,k}^{(n)}}{N},$$

which makes intuitive sense as the sample averages of our estimated probabilities for each possible value of the initial state. Next,

$$T_{k,l} = \frac{\sum_{n=1}^N \sum_{t=1}^{T-1} Q_{t,t+1,k,l}^{(n)}}{\sum_{n=1}^N \sum_{t=1}^{T-1} q_{tk}^{(n)}},$$

which has the intuitive interpretation of the (normalized) average of the transition probabilities, and finally

$$\pi_{k,m} = \frac{\sum_{n=1}^N \sum_{t=1}^T q_{t,k}^{(n)} x_{t,m}^{(n)}}{\sum_{n=1}^N \sum_{t=1}^T q_{t,k}^{(n)}},$$

which has the intuitive interpretation of a weighted average of the emissions given the state. After updating these parameters, we repeat the EM algorithm until convergence of said parameters.

## 1.6 Exercise: Parameter Estimation in Supervised HMMs

You are trying to predict the weather using an HMM. The hidden states are the weather of the day, which may be sunny or rainy, and the observable states are the color of the clouds, which can be white or gray. You have data on the weather and clouds from one sequence of four days (note: the hidden states are observed here):

Day	Weather	Clouds
1	Sunny	White
2	Rainy	Gray
3	Rainy	Gray
4	Sunny	Gray

1. Draw a graphical model representing the HMM.
2. Give the values of  $N, T, c$  and of the one-hot vectors  $\mathbf{z}_1^{(1)}, \dots, \mathbf{z}_4^{(1)}, \mathbf{x}_1^{(1)}, \dots, \mathbf{x}_4^{(1)}$ .
3. Estimate and interpret the values of the parameters  $\boldsymbol{\theta}, \mathbf{T}, \{\boldsymbol{\pi}\}$  using the MLE estimators for the supervised HMM provided in the previous subsection.

## 1.7 Exercise: EM for HMMs

You are trying to model a toy's state using an HMM. At each time step, the toy can be active (state 1) or inactive (state 2), but you can only observe the color of the indicator light, which can be red (observation state 1) or green (observation state 2). You have collected data from one sequence:

Time	Light
1	Green
2	Red
3	Green

You initialize your EM with  $\boldsymbol{\theta} = [\frac{1}{2} \ \frac{1}{2}]^\top$ ,  $\mathbf{T} = \begin{bmatrix} \frac{2}{3} & \frac{1}{3} \\ \frac{3}{3} & \frac{2}{3} \end{bmatrix}$ ,  $\boldsymbol{\pi}_1 = [\frac{1}{4} \ \frac{3}{4}]^\top$ ,  $\boldsymbol{\pi}_2 = [\frac{3}{4} \ \frac{1}{4}]^\top$ .

1. Compute  $\alpha_1, \alpha_2, \alpha_3, \beta_1, \beta_2, \beta_3$  for the forward-backward algorithm using the initial parameter values.

2. Refer to the definition of  $\mathbf{q}_t^{(1)}$  in Section 1.6.2. Now, compute the values of  $\mathbf{q}_1^{(1)}, \mathbf{q}_2^{(1)}$  using the  $\alpha$  and  $\beta$  values.
3. Refer to the definition of  $\mathbf{Q}_{t,t+1}^{(1)}$  in Section 1.6.2. Compute the value of  $\mathbf{Q}_{1,2}^{(1)}$  using the  $\alpha$  and  $\beta$  values.

During EM, at one point you obtain the following values after the E step:

$$\mathbf{q}_1^{(1)} = \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix}^\top, \quad \mathbf{q}_2^{(1)} = \begin{bmatrix} 1 & 2 \\ 3 & 3 \end{bmatrix}^\top, \quad \mathbf{q}_3^{(1)} = \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix}^\top$$

$$\mathbf{Q}_{1,2}^{(1)} = \begin{bmatrix} \frac{1}{6} & \frac{1}{2} \\ \frac{1}{6} & \frac{1}{6} \end{bmatrix}, \quad \mathbf{Q}_{2,3}^{(1)} = \begin{bmatrix} \frac{1}{6} & \frac{1}{6} \\ \frac{1}{2} & \frac{1}{6} \end{bmatrix}$$

4. Perform the M step by updating the parameters  $\boldsymbol{\theta}, \mathbf{T}, \boldsymbol{\pi}_1, \boldsymbol{\pi}_2$ .

## 2 Markov Decision Processes

A Markov Decision Process (MDP) is a framework for modeling an agent's actions in the world. It consists of:

1. A set of states  $S$
2. A set of actions  $A$
3. A reward function  $r : S \times A \rightarrow \mathbb{R}$
4. A transition model  $p(s'|s, a), \forall s, s' \in S, a \in A$ .

*Note:* Transitions to the next state only depend on the value of the current state (and the current action) and thus exhibit the *Markov Property*.

A *policy*  $\pi$  is a mapping from states to actions, i.e.  $\pi : S \rightarrow A$ .

### 2.1 Finite time horizon MDP

In the finite horizon setting, a policy may vary with the number of time periods remaining.  $\pi_{(t)}$  denotes the policy with  $t$  time steps to go.  $T$  is the decision horizon. The value of a policy with  $t$  time steps to go is defined inductively to be:

$$V_{(t)}^{\pi}(s) = \begin{cases} r(s, \pi_{(1)}(s)) & \text{if } t = 1 \\ r(s, \pi_{(t)}(s)) + \sum_{s' \in S} p(s'|s, \pi_{(t)}(s)) V_{(t-1)}^{\pi}(s') & \text{o.w.} \end{cases} \quad (1)$$

The process of computing these values inductively, working from the end of the horizon to the present, is called *value iteration*. If we instead look forward in time, we are computing the expected value of the policy

$$V_T^{\pi}(s) = \mathbb{E}_{s_1, \dots, s_T} \left[ \sum_{t=0}^T r(s_t, \pi_{(T-t)}(s_t)) \right] \quad (2)$$

by induction, where  $s_1 := s$ .  $V^{\pi}(s)$  is the MDP value function.

In an MDP, the general goal is to find an optimal policy by maximizing the expected reward under the policy, i.e. maximizing the value function. This is the *planning problem*.



## 2.2 Infinite Horizon MDP

**Policy Evaluation** We can also send  $T \rightarrow \infty$ , i.e. have an infinite time horizon. In that case, we need a discount factor  $0 < \gamma < 1$ , and we want to compute the value function

$$V^\pi(s) = \mathbb{E}_{s_1, s_2, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \right] \quad (3)$$

where  $s_1 := s$ , and the  $\gamma$  factor ensures convergence (assuming bounded rewards). In this setting, we only worry about stationary policies that don't vary with time. This is the *policy evaluation* problem; for any given policy  $\pi$ , we can find  $V^\pi(s)$  by solving the system of linear equations

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s' \in S} p(s'|s, \pi(s)) V^\pi(s') \quad (4)$$

These capture consistency about the value function. To solve this system, we can use Gaussian elimination, or simply iterate until convergence as in the finite horizon case.

Given a policy  $\pi$  and  $\theta$  (small positive number), we find  $V^\pi$  iteratively as follows:

- Initialize:  $V(s) = 0$  for all states  $s$ .
- Repeat
  - Update step:

$$V'(s) = r(s, \pi(s)) + \gamma \sum_{s' \in S} p(s'|s, \pi(s)) V(s'), \quad \forall s \quad (5)$$

- $\Delta = \max(|V' - V|)$
- $V \leftarrow V'$

until  $\Delta < \theta$

**Value Iteration** Suppose we have an optimal policy  $\pi^*$ . This satisfies the following set of equations known as the *Bellman equations*:

$$V^*(s) = \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^*(s') \right] \quad (6)$$

where  $V^* \triangleq V^{\pi^*}$ . Assuming we know  $V^*$ , we can read off the optimal policy by setting

$$\pi^*(s) = \arg \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^*(s') \right] \quad (7)$$

In order to find  $V^*$ , we can use *value iteration*:

- Initialize:  $V(s) = 0$  for all states  $s$ .
- Update step (Bellman operator):

$$V'(s) = \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V(s') \right], \quad \forall s \quad (8)$$

- $V \leftarrow V'$

where we iterate until convergence of  $V$ , which is guaranteed. With our converged  $V$ , we can then find  $\pi^*$  as in Equation 7.

**Policy Iteration** Another approach to planning is called *policy iteration*. To do policy iteration, we evaluate a proposed policy  $\pi$  by finding  $V^\pi$  as in Equation 4. This is Evaluation step (E step). Then, we do a policy improvement step (I step) by the equation

$$\pi'(s) \leftarrow \arg \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^\pi(s') \right], \quad \forall s \quad (9)$$

We repeat the E and I steps until the policy  $\pi$  converges (stops changing).

**Comparing Value and Policy Iteration** Policy iteration involves two distinct steps: policy evaluation and policy improvement, applied iteratively to refine the policy  $\pi$ . Value iteration combines these steps by performing a simplified update of the value function using the Bellman optimality operator. This update implicitly improves the policy as the value function approaches optimality. Policy iteration takes more computation per iteration, but tends to converge faster in practice.

## 2.3 Exercise: Markov Decision Process

(Sutton & Barto 2012) Consider an MDP on the following grid:

	A	
	B	

At each square, we can go left, right, up, or down. Normally we get a reward of 0 from moving, but if we attempt to move off the grid, we get a reward of  $-1$  and stay where we are. Also, if we move onto square A, we get a reward of 10 and are teleported to square B.

Suppose our actions also fail with probability 0.5, i.e. with probability 0.5 we stay on the current square. Also suppose our MDP is infinite horizon, and take  $\gamma = 0.9$  to be the discount factor.

1. **Defining the MDP** Identify the states  $S$ , actions  $A$ , rewards, and transition probabilities  $p(s'|s, a)$  in this problem.
2. **Policy Evaluation** Suppose  $\pi$  is the policy where we always choose to go right. Write the equations to find the values  $V^\pi(s)$ .
3. **Value Iteration** Write the second iteration of value iteration, i.e. starting by initializing  $V(s) = \max_{a \in A} r(s, a)$ .
4. **Policy Iteration** Write the first iteration of policy iteration, starting with  $V^\pi(s) = 0$  for all  $s$ . (We could also initialize a policy, and do the Evaluation step to get started.)

### 3 Kalman Filters (Bonus Material)

Now consider the following dynamical system model:

$$z_{t+1} = \Phi z_t + \epsilon_t$$

$$x_t = Az_t + \gamma_t$$

where  $z$  are the hidden variables and  $x$  are the observed measurements.  $\Phi$  and  $A$  are known constants, while  $\epsilon$  and  $\gamma$  are random variables drawn from the following normal distributions:

$$\epsilon_t \sim \mathcal{N}(\mu_\epsilon, \sigma_\epsilon^2)$$

$$\gamma_t \sim \mathcal{N}(\mu_\gamma, \sigma_\gamma^2)$$

This is called a (one-dimensional) linear Gaussian state-space model. It is closely related to an HMM – try drawing out the graphical model! – but here the hidden states and the observations are now continuous and normally distributed. Linear Gaussian state-space models have convenient mathematical properties and can be used to describe noisy measurements of a moving object (e.g. missiles, rodents, hands), market fluctuations, etc.

The Kalman filter is an algorithm to perform filtering in linear Gaussian state-space models, i.e. to find the distribution of  $z_t$  given observations  $x_1, \dots, x_t$ . The distribution of  $z_t | x_1, \dots, x_s$  will be  $\mathcal{N}(\mu_{t|s}, \sigma_{t|s}^2)$ . If we start with  $\mu_{t-1|t-1}$  and  $\sigma_{t-1|t-1}^2$ , the algorithm tells us to

1. Define the distribution of  $z_t | x_1, \dots, x_{t-1}$  by computing  $\mu_{t|t-1}$  and  $\sigma_{t|t-1}^2$ . This is called the prediction step.
2. Define the distribution of  $z_t | x_1, \dots, x_t$  by computing  $\mu_{t|t}$  and  $\sigma_{t|t}^2$ . This is called the update step.

The Kalman filter alternates between prediction and update steps, assimilating observations one at a time. It requires one forward pass through the data, and is analogous to obtaining the  $\alpha$ 's in an HMM.