

CS 181 Spring 2021 Section 11

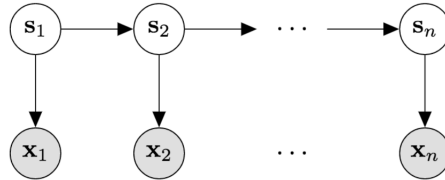
HMMs, Kalman Filters, and MDP's

Solution

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



Consider a sequence of one-hot encoded states $\mathbf{s}_1, \dots, \mathbf{s}_n$ where $\mathbf{s}_t \in \{S_k\}_{k=1}^c$, and a corresponding sequence of observations $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ where $\mathbf{x}_t \in \{O_j\}_{j=1}^m$. Each state can be one of c possible states, and each observation can be one of m possible observations. Note that N is the number of data points (each of which is a sequence), where n is the length of a sequence (assume all sequences are the same length).

1.2 Model Assumptions

HMMs are characterized by and allow us to reason about the following joint distribution

$$p(\mathbf{s}_1, \dots, \mathbf{s}_n, \mathbf{x}_1, \dots, \mathbf{x}_n) = p(\mathbf{s}_1, \dots, \mathbf{s}_n)p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{s}_1, \dots, \mathbf{s}_n)$$

However, it's not immediately obvious how we should optimize this model, and the following assumptions make this easier:

- The future hidden state is independent of past hidden states given the present (Markov Property):

$$p(\mathbf{s}_{t+1} | \mathbf{s}_1, \dots, \mathbf{s}_t, \mathbf{x}_1, \dots, \mathbf{x}_t) = p(\mathbf{s}_{t+1} | \mathbf{s}_t)$$

- Observations only depend on the present hidden state:

$$p(\mathbf{x}_t | \mathbf{s}_1, \dots, \mathbf{s}_t, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) = p(\mathbf{x}_t | \mathbf{s}_t)$$

Notice that the above assumptions allow us to factor the joint as follows:

$$\begin{aligned} p(\mathbf{s}_1, \dots, \mathbf{s}_n, \mathbf{x}_1, \dots, \mathbf{x}_n) &= \\ p(\mathbf{s}_1, \dots, \mathbf{s}_n) p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{s}_1, \dots, \mathbf{s}_n) &= \\ p(\mathbf{s}_1) \prod_{t=1}^{n-1} p(\mathbf{s}_{t+1} | \mathbf{s}_t) \prod_{t=1}^n p(\mathbf{x}_t | \mathbf{s}_t) & \end{aligned}$$

1.3 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

Solution

1. Stock market price data: Yes, an HMM is appropriate since stock market data is time-dependent. Observed data: stock prices listed on exchanges. Hidden states: true value of the stock, perhaps a combination of company policies, growth potential, economic conditions, etc.
2. Recommendations on a database of movie reviews: No, an HMM would not be appropriate since we don't expect user preferences to change much over time.
3. Daily precipitation data in Boston: Yes, precipitation today is very likely to affect the chance of precipitation tomorrow. Observed data: amount of precipitation each day. Possible hidden states: true weather conditions, such as humidity or chance of rain.
4. Optical character recognition, where we are identifying words: Yes, word recognition is very dependent upon the sequence of characters. Observed data: image pixels of written characters. Hidden states: the true character represented (think MNIST from the last theory pset).

End Solution

-3

1.4 Parameterization

- $\theta \in \mathbb{R}^c$: defines the prior distribution over initial hidden states
- $\mathbf{T} \in \mathbb{R}^{c \times c}$: transition matrix where T_{kj} is the probability of transitioning from S_k to S_j
- $\{\boldsymbol{\pi}\}_{k=1}^c$: conditional probabilities of observations given hidden states such that $p(\mathbf{x}_t = O_j | \mathbf{s}_t = S_k; \{\boldsymbol{\pi}\}) = \pi_{kj}$. $\forall k \boldsymbol{\pi}_k \in \mathbb{R}^m$.

First, we need to estimate the parameters from the data, which we can do with a variant of EM. Then, with our trained HMM, we are able to perform several inference tasks on our data.

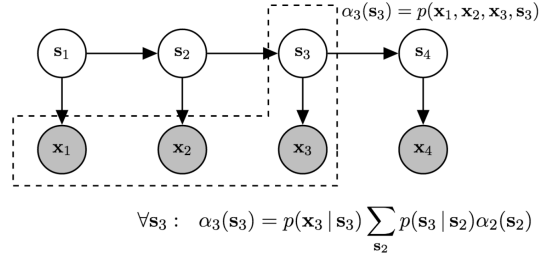
1.5 EM for HMMs

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

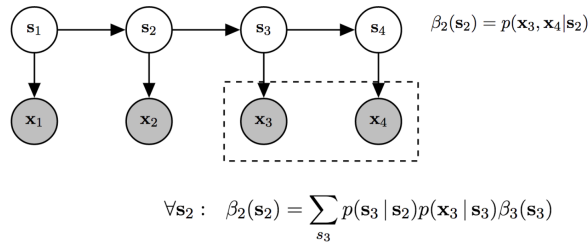
1.5.1 Forward-Backward Algorithm

The HMM model is characterized by the joint distribution $p(\mathbf{s}_1, \dots, \mathbf{s}_n, \mathbf{x}_1, \dots, \mathbf{x}_n)$, 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. We define the recurrence relations $\alpha_t(\mathbf{s}_t)$ and $\beta_t(\mathbf{s}_t)$ in the E-Step:

- $\alpha_t(\mathbf{s}_t)$ represents the joint probability of observations $1, \dots, t$ and state t . α_t can be defined in terms of α_{t-1} . We move **forwards** through the sequence to calculate the α 's
- $\beta_t(\mathbf{s}_t)$ represents the joint probability of observations $t+1, \dots, n$ conditioned on state t . β_t can be defined in terms of β_{t+1} . We move **backwards** through the sequence to calculate the β 's.



(a) alpha



(b) beta

Note that the probabilities we use for calculating α and β are given by the parameters that we fix in the E-Step.

$$\forall \mathbf{s}_t : \alpha_t(\mathbf{s}_t) = \begin{cases} p(\mathbf{x}_t | \mathbf{s}_t) \sum_{\mathbf{s}_{t-1}} p(\mathbf{s}_t | \mathbf{s}_{t-1}) \alpha_{t-1}(\mathbf{s}_{t-1}) & \text{if } 1 < t \leq n \\ p(\mathbf{x}_1 | \mathbf{s}_1) p(\mathbf{s}_1) & \text{o.w.} \end{cases}$$

$$\forall \mathbf{s}_t : \beta_t(\mathbf{s}_t) = \begin{cases} \sum_{\mathbf{s}_{t+1}} p(\mathbf{s}_{t+1} | \mathbf{s}_t) p(\mathbf{x}_{t+1} | \mathbf{s}_{t+1}) \beta_{t+1}(\mathbf{s}_{t+1}) & \text{if } 1 \leq t < n \\ 1 & \text{o.w.} \end{cases}$$

1.5.2 Inference Patterns with α, β

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

- $\alpha_t(\mathbf{s}_t) \beta_t(\mathbf{s}_t) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}_t) \propto p(\mathbf{s}_t | \mathbf{x}_1, \dots, \mathbf{x}_n)$
- joint of observations: $p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{\mathbf{s}_t} \alpha_t(\mathbf{s}_t) \beta_t(\mathbf{s}_t)$ (for any t)
- smoothing: $p(\mathbf{s}_t | \mathbf{x}_1, \dots, \mathbf{x}_n) \propto p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}_t) = \alpha_t(\mathbf{s}_t) \beta_t(\mathbf{s}_t)$
- prediction: $p(\mathbf{x}_{n+1} | \mathbf{x}_1, \dots, \mathbf{x}_n) \propto \sum_{\mathbf{s}_n, \mathbf{s}_{n+1}} \alpha_n(\mathbf{s}_n) p(\mathbf{s}_{n+1} | \mathbf{s}_n) p(\mathbf{x}_{n+1} | \mathbf{s}_{n+1})$
- transition: $p(\mathbf{s}_t, \mathbf{s}_{t+1} | \mathbf{x}_1, \dots, \mathbf{x}_n) \propto \alpha_t(\mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t) p(\mathbf{x}_{t+1} | \mathbf{s}_{t+1}) \beta_{t+1}(\mathbf{s}_{t+1})$

1.5.3 E-Step

The goal of the expectation step is to compute the expected values of the hidden states given a fixed set of parameters $\mathbf{w} = \{\mathbf{T}, \boldsymbol{\theta}, \{\boldsymbol{\pi}_k\}\}$. That is, we estimate the state distribution for $\mathbf{s}_1^i, \dots, \mathbf{s}_n^i$ given \mathbf{x}^i .

Let the $c \times 1$ vector $\mathbf{q}_t^i = (q_{t1}^i, \dots, q_{tc}^i)$ represent \mathbf{x}^i 's distribution over states for time t under the current parameters. Let $\mathbf{Q}_{t,t+1}^i$ be the $c \times c$ matrix of transition probabilities under the current parameters. Then

- α 's and β 's are defined in terms of fixed parameters.
- \mathbf{q} 's are defined in terms of α 's and β 's
- Calculate $q_{tk}^i = p(\mathbf{s}_t^i = S_k | \mathbf{x}^i; \mathbf{w})$ for all t and k (use smoothing eq. just above)
- Calculate $q_{t,t+1,k,\ell}^i = p(\mathbf{s}_t^i = S_k, \mathbf{s}_{t+1}^i = S_\ell | \mathbf{x}^i; \mathbf{w})$ (use transition eq. just above)

1.5.4 M-Step

Now we need to update our parameters to maximize the expected complete-data log likelihood $\mathbb{E}_{\mathbf{S}}[\ln p(\mathbf{x}, \mathbf{S}; \mathbf{w})]$. Applying the appropriate Lagrange multipliers and maximizing with respect to each of the parameters of interest, we recover the following update equations:

$$\hat{N}_{1k} = \sum_{i=1}^N q_{1k}^i \text{ (first period)} \quad \text{and more generally} \quad \hat{N}_k = \sum_{i=1}^N \sum_{t=1}^n q_{tk}^i \text{ (all periods)}$$

$$\hat{N}_{-nk} = \sum_{i=1}^N \sum_{t=1}^{n-1} q_{tk}^i \text{ (without last period)}$$

$$\hat{N}_{k\ell} = \sum_{i=1}^N \sum_{t=1}^{n-1} q_{t,t+1,k,\ell}^i \text{ (transitions)}$$

$$\hat{N}_{kj} = \sum_{i=1}^N \sum_{t=1}^n q_{tk}^i x_{tj}^i \text{ (observations)}$$

$$\hat{\theta}_k = \frac{\hat{N}_{1k}}{N} \quad \hat{\pi}_{kj} = \frac{\hat{N}_{kj}}{\hat{N}_k} \quad \hat{t}_{k\ell} = \frac{\hat{N}_{k\ell}}{\hat{N}_{-nk}}$$

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, n, c and of the one-hot vectors $\mathbf{s}_1^1, \dots, \mathbf{s}_4^1, \mathbf{x}_1^1, \dots, \mathbf{x}_4^1$.
3. Estimate and interpret the values of the parameters $\boldsymbol{\theta}, \mathbf{T}, \{\boldsymbol{\pi}_k\}_{k=1}^c$ using the MLE estimators for the supervised HMM:

$$\hat{\theta}_k = \frac{N_{1k}}{N}, \quad \hat{t}_{kl} = \frac{N_{kl}}{N_{-nk}}, \quad \hat{\pi}_{kj} = \frac{N_{kj}}{N_k}$$

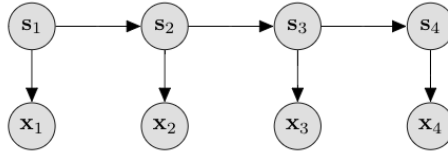
$$N_k = \sum_{i=1}^N \sum_{t=1}^n s_{tk}^i, \quad N_{1k} = \sum_{i=1}^N s_{1,k}^i, \quad N_{-nk} = \sum_{i=1}^N \sum_{t=1}^{n-1} s_{tk}^i$$

$$N_{kl} = \sum_{i=1}^N \sum_{t=1}^{n-1} s_{t,k}^i s_{t+1,l}^i, \quad N_{kj} = \sum_{i=1}^N \sum_{t=1}^n s_{tk}^i x_{tj}^i$$

3

Solution

1.



(all nodes are observed)

2. $N = 1$, the number of sequences observed

$n = 4$, the length of the sequences

$c = 2$, the number of states a hidden state can take

$\mathbf{s}_1^1 = [1 \ 0]^\top, \mathbf{s}_2^1 = [0 \ 1]^\top, \mathbf{s}_3^1 = [0 \ 1]^\top, \mathbf{s}_4^1 = [1 \ 0]^\top$

$\mathbf{x}_1^1 = [1 \ 0]^\top, \mathbf{x}_2^1 = [0 \ 1]^\top, \mathbf{x}_3^1 = [0 \ 1]^\top, \mathbf{x}_4^1 = [0 \ 1]^\top$

3. $\hat{\boldsymbol{\theta}} = [1 \ 0]^\top$, the distribution of the weather for the initial state

$\hat{\mathbf{T}} = \begin{bmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$, the transition probabilities for the weather

$\hat{\boldsymbol{\pi}}_1 = [\frac{1}{2} \ \frac{1}{2}]^\top$, the distribution of cloud colors on sunny days

$\hat{\boldsymbol{\pi}}_2 = [0 \ 1]^\top$, the distribution of cloud colors on rainy days

End Solution

-3

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 $\theta = [\frac{1}{2} \ \frac{1}{2}]^\top$, $\mathbf{T} = \begin{bmatrix} \frac{2}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{2}{3} \end{bmatrix}$, $\pi_1 = [\frac{1}{4} \ \frac{3}{4}]^\top$, $\pi_2 = [\frac{3}{4} \ \frac{1}{4}]^\top$.

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. How is \mathbf{q}_t^1 defined? Compute the values of $\mathbf{q}_1^1, \mathbf{q}_2^1$ using the α and β values.
3. How is $\mathbf{Q}_{t,t+1}^1$ defined? Compute the value of $\mathbf{Q}_{1,2}^1$ using the α and β values.

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

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

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

1. Use the above values to compute $\hat{N}_k, \hat{N}_{kl}, \hat{N}_{kj}$.
2. Complete the M step by updating the parameters $\theta, \mathbf{T}, \pi_1, \pi_2$.

Solution

1. Using the recursive definitions for α, β and the current values of $\boldsymbol{\theta}, \mathbf{T}, \boldsymbol{\pi}_1, \boldsymbol{\pi}_2$:

$$\begin{aligned}\alpha_1(\mathbf{s}_1^1) &= \begin{cases} \frac{3}{8} & \mathbf{s}_1^1 = \text{active} \\ \frac{1}{8} & \mathbf{s}_1^1 = \text{inactive} \end{cases} \\ \alpha_2(\mathbf{s}_2^1) &= \begin{cases} \frac{7}{96} & \mathbf{s}_2^1 = \text{active} \\ \frac{15}{96} & \mathbf{s}_2^1 = \text{inactive} \end{cases} \\ \alpha_3(\mathbf{s}_3^1) &= \begin{cases} \frac{29}{384} & \mathbf{s}_3^1 = \text{active} \\ \frac{37}{1152} & \mathbf{s}_3^1 = \text{inactive} \end{cases} \\ \beta_3(\mathbf{s}_3^1) &= \begin{cases} 1 & \mathbf{s}_3^1 = \text{active} \\ 1 & \mathbf{s}_3^1 = \text{inactive} \end{cases} \\ \beta_2(\mathbf{s}_2^1) &= \begin{cases} \frac{7}{12} & \mathbf{s}_2^1 = \text{active} \\ \frac{5}{12} & \mathbf{s}_2^1 = \text{inactive} \end{cases} \\ \beta_1(\mathbf{s}_1^1) &= \begin{cases} \frac{29}{144} & \mathbf{s}_1^1 = \text{active} \\ \frac{37}{144} & \mathbf{s}_1^1 = \text{inactive} \end{cases}\end{aligned}$$

2. q_{tk}^1 is the probability that \mathbf{s}_t^1 is S_k (given the observations), and $q_{tk}^1 = p(\mathbf{s}_t^1 = S_k \mid \mathbf{x}^1; \mathbf{w}) \propto \alpha_t(S_k)\beta_t(S_k)$. Then $\mathbf{q}_1^1 \propto [\frac{87}{1152}, \frac{37}{1152}]^\top$, so $\mathbf{q}_1^1 = [\frac{87}{124}, \frac{37}{124}]^\top$. Also, $\mathbf{q}_2^1 \propto [\frac{49}{1152}, \frac{75}{1152}]^\top$, so $\mathbf{q}_2^1 = [\frac{49}{124}, \frac{75}{124}]^\top$.
3. $q_{t,t+1,k,l}^1$ is the probability that \mathbf{s}_t^1 is S_k and \mathbf{s}_{t+1}^1 is S_l (given the observations), and $q_{t,t+1,k,l}^1 = p(\mathbf{s}_t^1 = S_k, \mathbf{s}_{t+1}^1 = S_l \mid \mathbf{x}^1; \mathbf{w}) \propto \alpha_t(\mathbf{s}_t)p(\mathbf{s}_{t+1} \mid \mathbf{s}_t)p(\mathbf{x}_{t+1} \mid \mathbf{s}_{t+1})\beta_{t+1}(\mathbf{s}_{t+1})$. Then

$$\mathbf{Q}_{1,2}^1 \propto \begin{bmatrix} \frac{42}{1152} & \frac{45}{1152} \\ \frac{7}{1152} & \frac{30}{1152} \end{bmatrix}$$

so

$$\mathbf{Q}_{1,2}^1 = \begin{bmatrix} \frac{42}{124} & \frac{45}{124} \\ \frac{7}{124} & \frac{30}{124} \end{bmatrix}$$

- 4.

$$\text{For } \hat{N}_k: \quad \hat{N}_1 = \frac{5}{3}, \quad \hat{N}_2 = \frac{4}{3}$$

$$\text{For } \hat{N}_{kl}: \quad \hat{N}_{1,1} = \frac{1}{3}, \quad \hat{N}_{1,2} = \frac{2}{3}, \quad \hat{N}_{2,1} = \frac{2}{3}, \quad \hat{N}_{2,2} = \frac{1}{3}$$

$$\text{For } \hat{N}_{kj}: \quad \hat{N}_{1,1} = \frac{1}{3}, \quad \hat{N}_{1,2} = \frac{4}{3}, \quad \hat{N}_{2,1} = \frac{2}{3}, \quad \hat{N}_{2,2} = \frac{2}{3}$$

$$\begin{aligned}
 5. \quad \boldsymbol{\theta} &= \begin{bmatrix} 2 \\ 3 \\ 3 \end{bmatrix}^T \\
 \mathbf{T} &= \begin{bmatrix} 1 & 2 \\ 3 & 3 \\ 2 & 1 \\ 3 & 3 \end{bmatrix} \\
 \boldsymbol{\pi}_1 &= \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix}^T \\
 \boldsymbol{\pi}_2 &= \begin{bmatrix} 4 \\ 1 \\ 2 \end{bmatrix}^T
 \end{aligned}$$

End Solution

-3

2 Kalman Filters

Now consider the following dynamical system model:

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

$$x_t = A z_t + \gamma_t$$

where z are the hidden variables and x are the observed measurements. Φ and A are known constants, while ϵ and γ 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 α 's in an HMM.

3 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 $p(s'|s, a), \forall s, s' \in S, a \in A$.

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

3.1 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

$$V^\pi(s) = \mathbb{E}_{s_1, \dots, s_T} \left[\sum_{t=1}^T r(s_t, \pi_{(T+1-t)}(s_t)) \right] \quad (2)$$

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.

3.2 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

$$V^\pi(s) = \mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, \pi(s_t)) \right] \quad (3)$$

where $s_1 := s$, and the γ 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 π , 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 π and θ (small positive number), we find V^π 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 π^* . 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 π^* as in Equation 7.

Policy Iteration Another approach to planning is called *policy iteration*. To do policy iteration, we evaluate a proposed policy π by finding V^π 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 π converges (stops changing).

Note that policy iteration takes more computation per iteration, but tends to converge faster in practice.

3.3 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 π 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.)

Solution:

1. Defining the MDP

There are 9 states, one for each square of the grid, which we will denote as (i, j) for $i = 0, 1, 2$ (rows) and $j = 0, 1, 2$ (columns), where $(0, 0)$ is the top left corner. There are 4 actions (one for each direction), call them L,R,U,D.

We will define $r(s, a)$ for a state-action pair as the expected reward from taking action a in state s ; for example, $r((0, 0), R) = 0.5 \times 0 + 0.5 \times 10 = 5$.¹

The probabilities of transition between any two valid states are all 0.5, and the probabilities of remaining at any state are all 0.5 (except for moving off the edge of the board, where the probability is 1 of staying). In addition, a successful transition into state A results in the agent being in state B (teleportation!).

2. Policy Evaluation

Recall the policy evaluation equations

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

We have

$$\begin{aligned} V^\pi((0, 0)) &= 5 + (0.9)(0.5 \times V^\pi((0, 0)) + 0.5 \times V^\pi((2, 1))) \\ V^\pi((0, 1)) &= 0 + (0.9)(0.5 \times V^\pi((0, 1)) + 0.5 \times V^\pi((0, 2))) \\ V^\pi((0, 2)) &= -0.5 + (0.9)V^\pi((0, 2)) \\ V^\pi((1, 0)) &= 0 + (0.9)(0.5 \times V^\pi((1, 0)) + 0.5 \times V^\pi((1, 1))) \\ V^\pi((1, 1)) &= 0 + (0.9)(0.5 \times V^\pi((1, 1)) + 0.5 \times V^\pi((1, 2))) \\ V^\pi((1, 2)) &= -0.5 + (0.9)V^\pi((1, 2)) \\ V^\pi((2, 0)) &= 0 + (0.9)(0.5 \times V^\pi((2, 0)) + 0.5 \times V^\pi((2, 1))) \\ V^\pi((2, 1)) &= 0 + (0.9)(0.5 \times V^\pi((2, 1)) + 0.5 \times V^\pi((2, 2))) \\ V^\pi((2, 2)) &= -0.5 + (0.9)V^\pi((2, 2)) \end{aligned}$$

We can compactly write this in matrix form as follows:

$$\mathbf{V} = \mathbf{r} + \gamma \mathbf{T} \mathbf{V}$$

where

$$\begin{aligned} \mathbf{V} = [&V^\pi((0, 0)), V^\pi((0, 1)), V^\pi((0, 2)), \\ &V^\pi((1, 0)), V^\pi((1, 1)), V^\pi((1, 2)), \\ &V^\pi((2, 0)), V^\pi((2, 1)), V^\pi((2, 2))]^\top \end{aligned}$$

¹An earlier version of this handout treated the reward as if the robot could control its environment, but its actuator is noisy. Thus, we should adopt $r(s, a)$ to be defined as the expected reward for taking action a in state s .

$$\mathbf{r} = [5, 0, -0.5, 0, 0, -0.5, 0, 0, -0.5]^\top$$

$$\mathbf{T} = \begin{bmatrix} 0.5 & 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0 \\ 0 & 0.5 & 0.5 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 0.5 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.5 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0.5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Note that the policy affects \mathbf{r} and \mathbf{T} . This linear system can be solved using standard methods (e.g. Gaussian elimination). We find that (approximately)

$$\mathbf{V} = [5.7, -4.1, -5, -3.4, -4.1, -5, -3.4, -4.1, -5]$$

3. Value Iteration

Note that initially we have

$$\begin{aligned} V((0,0)) &= 5, & V((0,1)) &= 0, & V((0,2)) &= 5 \\ V((1,0)) &= 0, & V((1,1)) &= 5, & V((1,2)) &= 0 \\ V((2,0)) &= 0, & V((2,1)) &= 0, & V((2,2)) &= 0 \end{aligned}$$

For state $(0,0)$ we have the update:

$$\begin{aligned} V'((0,0)) &= \max[-0.5 + (0.9) \times V((0,0)), \\ &\quad -0.5 + (0.9) \times V((0,0)), \\ &\quad 5 + (0.9)(0.5 \times V((0,0)) + 0.5 \times V((2,1))), \\ &\quad 0 + (0.9)(0.5 \times V((0,0)) + 0.5 \times V((1,0)))] \\ &= \max[-0.5 + (0.9)(5), \\ &\quad -0.5 + (0.9)(5), \\ &\quad 5 + (0.9)(0.5 \times 5 + 0.5 \times 0), \\ &\quad 0 + (0.9)(0.5 \times 5 + 0.5 \times 0)] \\ &= 7.25 \end{aligned}$$

by taking the max over each of the four directions (left, up, right, down in that order). You can verify that the following hold for the

remaining states:

$$\begin{aligned} V'((0, 1)) &= 2.25, V'((0, 2)) = 7.25 \\ V'((1, 0)) &= 2.25, V'((1, 1)) = 7.25, V'((1, 2)) = 2.25 \\ V'((2, 0)) &= 0, V'((2, 1)) = 2.25, V'((2, 2)) = 0 \end{aligned}$$

4. Policy Iteration

Let's first do the I step (policy Improvement) to find the new policy π' . For $(0, 0)$, considering that we assume $V^\pi(s)$ under the current policy, we have

$$\pi'((0, 0)) = \arg \max_{a \in A} [-0.5, -0.5, 5, 0] = R$$

For other states, we get

$$\begin{aligned} \pi'((0, 0)) &= R, \pi'((0, 1)) = R, \pi'((0, 2)) = L \\ \pi'((1, 0)) &= R, \pi'((1, 1)) = U, \pi'((1, 2)) = D \\ \pi'((2, 0)) &= R, \pi'((2, 1)) = R, \pi'((2, 2)) = U \end{aligned}$$

breaking ties arbitrarily.

Next we do the E step (Evaluation). This is policy evaluation, and we can write our equation as $\mathbf{V} = \mathbf{r} + \gamma \mathbf{T} \mathbf{V}$, where based on the current policy π' we have

$$\begin{aligned} \mathbf{r} &= [5, 0, 5, 0, 5, 0, 0, 0, 0]^\top \\ \mathbf{T} &= \begin{bmatrix} 0.5 & 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0 \\ 0 & 0.5 & 0.5 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.5 & 0 & 0 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 0.5 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.5 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.5 & 0 & 0 & 0.5 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0.5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0 & 0 & 0 & 0.5 & 0 & 0 & 0.5 \end{bmatrix} \end{aligned}$$

Solving for \mathbf{V} , we find that

$$\mathbf{V} = [9.1, 7.5, 9.1, 7.5, 9.1, 0, 0, 0, 0]^\top$$

is the new estimate for $V^{\pi'}$.