

# Homework 5 - Introduction to Machine Learning For Engineers

kipngeno koech - bkoech

April 28, 2025

## 1 Gaussian Mixture Models

Consider an exponential mixture model for a 1-D dataset  $\{x_n\}$  with the density function

$$p(x) = \sum_{k=1}^K \omega_k \text{Exp}(x|\mu_k),$$

where  $K$  is the number of mixture components,  $\mu_k$  is the rate parameter, and  $\omega_k$  is the mixture weight corresponding to the  $k$ -th component. The exponential distribution is given by

$$\text{Exp}(x|\mu) = \mu \exp(-x\mu) \quad \text{for all } x \geq 0. \quad (1)$$

We would like to derive the model parameters  $(\omega_k, \mu_k)$  for all  $k$  using the EM algorithm. Consider the hidden labels  $z_n \in \{1, \dots, K\}$  and indicator variables  $r_{nk}$  that are 1 if  $z_n = k$  and 0 otherwise. The complete log-likelihood (assuming base  $e$  for the log) is then written as

$$\sum_n \log p(x_n, z_n) = \sum_n \sum_{k=1}^K [\log p(z_n = k) + \log p(x_n | z_n = k)].$$

1. Write down and simplify the expression for the complete log-likelihood for the exponential mixture model described above. Write and simplify the complete log-likelihood for the exponential mixture model.

The joint probability is:

$$p(x_n, z_n = k) = p(z_n = k)p(x_n | z_n = k) = \omega_k \text{Exp}(x_n | \mu_k) = \omega_k \mu_k \exp(-\mu_k x_n).$$

The log-likelihood for a single data point is:

$$\log p(x_n, z_n = k) = \log \omega_k + \log(\mu_k \exp(-\mu_k x_n)) = \log \omega_k + \log \mu_k - \mu_k x_n.$$

Using indicator variables  $r_{nk}$ , the complete log-likelihood over all data points is:

$$\ell = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log p(x_n, z_n = k) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} [\log \omega_k + \log \mu_k - \mu_k x_n].$$

2. Solve the M step of the EM algorithm and find  $\mu_k$  for  $k = 1, \dots, K$  that maximizes the complete log-likelihood. Taking the derivative of the log-likelihood with respect to  $\mu_k$  and setting it to zero, we have: Maximize the complete log-likelihood with respect to  $\mu_k$  in the M-step, subject to  $\sum_{k=1}^K \omega_k = 1$ .

The complete log-likelihood is:

$$\ell = \sum_{n=1}^N \sum_{k=1}^K r_{nk} [\log \omega_k + \log \mu_k - \mu_k x_n].$$

Differentiate with respect to  $\mu_k$ :

$$\frac{\partial \ell}{\partial \mu_k} = \sum_{n=1}^N r_{nk} \left[ \frac{1}{\mu_k} - x_n \right].$$

Set the derivative to zero:

$$\sum_{n=1}^N r_{nk} \frac{1}{\mu_k} = \sum_{n=1}^N r_{nk} x_n.$$

Solve for  $\mu_k$ :

$$\mu_k = \frac{\sum_{n=1}^N r_{nk}}{\sum_{n=1}^N r_{nk} x_n}.$$

For  $\omega_k$ , maximize  $\ell$  subject to  $\sum_{k=1}^K \omega_k = 1$  using a Lagrange multiplier:

$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log \omega_k + \lambda \left( 1 - \sum_{k=1}^K \omega_k \right).$$

Differentiate with respect to  $\omega_k$ :

$$\frac{\partial \mathcal{L}}{\partial \omega_k} = \sum_{n=1}^N \frac{r_{nk}}{\omega_k} - \lambda = 0 \implies \sum_{n=1}^N r_{nk} = \lambda \omega_k.$$

Sum over  $k$ :

$$\sum_{k=1}^K \sum_{n=1}^N r_{nk} = \lambda \sum_{k=1}^K \omega_k = \lambda.$$

Since  $\sum_{n=1}^N \sum_{k=1}^K r_{nk} = N$  and  $\sum_{k=1}^K \omega_k = 1$ , we get  $\lambda = N$ . Thus:

$$\omega_k = \frac{\sum_{n=1}^N r_{nk}}{\sum_{n=1}^N r_{nk} x_n}$$

3. Perform the E step of the EM algorithm and write the equation to update the soft labels  $r_{nk} = P(z_n = k | x_n)$ . Using Bayes' rule, we have: Compute the soft labels  $r_{nk} = P(z_n = k | x_n)$  in the E-step.

Using Bayes' rule:

$$r_{nk} = P(z_n = k | x_n) = \frac{P(x_n, z_n = k)}{P(x_n)} = \frac{p(z_n = k)p(x_n | z_n = k)}{\sum_{j=1}^K p(z_n = j)p(x_n | z_n = j)}.$$

Substitute:

$$p(z_n = k) = \omega_k, \quad p(x_n | z_n = k) = \mu_k \exp(-\mu_k x_n).$$

Thus:

$$r_{nk} = \frac{\omega_k \mu_k \exp(-\mu_k x_n)}{\sum_{j=1}^K \omega_j \mu_j \exp(-\mu_j x_n)}.$$

## 2 Eigen Faces

Face recognition is an important task in computer vision and machine learning. In this question, you will implement a classical approach called Eigenfaces. You will use face images from the Yale Face Database B, which contains face images from 10 people under 64 lighting conditions. Please include your code in the final PDF you turn in for full credit.

1. **Dataset.** Download the data file `face_data.mat`. It contains three sets of variables:
  - **image:** each element is a face image ( $50 \times 50$  matrix). You can use `matplotlib.pyplot.imshow` to visualize the image. The data is stored in a cell array.
  - **personID:** each element is the ID of the person, which takes values from 1 to 10.
  - **subsetID:** each element is the ID of the subset which takes values from 1 to 5. Here the face images are divided into 5 subsets. Each subset contains face images from all people, but with different lighting conditions.
2. **[10 points]** Implement PCA. Fill in the function `pca_fun` in the `pca.py` file. The function takes the data matrix (each row being a sample) and target dimensionality  $d$  (lower than or equal to the original dimensionality) as the input, and outputs the selected eigenvectors.
3. **[25 points]** Compute Eigenfaces. Take each  $50 \times 50$  training image and vectorize it into a 2500- dimensional vector. Use your PCA implementation from part (b) to perform PCA on all vectorized face images, and retain the first  $d = 200$  eigenvectors. These eigenvectors are called eigenfaces (when displayed as images). Please display the top 5 eigenfaces (use `imshow`) in your report.

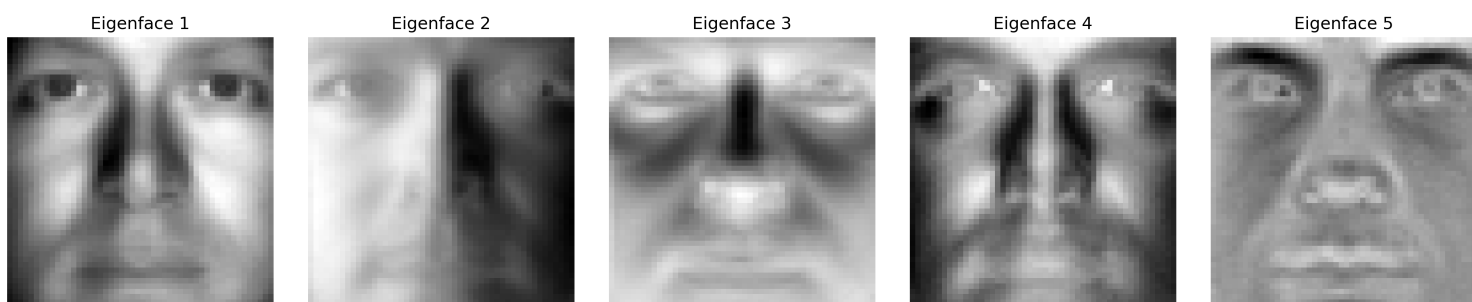


Figure 1: Top 5 Eigenfaces

### 3 Thompson Sampling

Consider the Thompson Sampling (TS) algorithm, a Bayesian approach to the multi-armed bandit problem. Consider a Bernoulli bandit with  $n$  arms, where each arm  $i$  at time-step  $1 \leq t \leq T$  has Bernoulli i.i.d. rewards  $r_{i,t} \in \{0, 1\}$  with  $\mathbb{E}[r_{i,t}] = \mu_i$ .

The TS algorithm starts with a prior distribution of  $\mu_i$  for each arm  $i$  using the  $P_{i,0} \sim \text{Beta}(1, 1)$  distribution and proceeds by selecting an arm based on the posterior distribution as follows. Note the prior distribution of  $\mu_i$  at time  $t$  is denoted as  $P_{i,t-1}$  and the posterior as  $P_{i,t}$ . Further, the posterior of the current time-step becomes the prior for the next time-step.

---

**Algorithm 1** Thompson Sampling

---

```

1: for  $t = 1, 2, \dots, T$  do
2:   Sample  $\hat{\mu}_{i,t} \sim P_{i,t-1}$  for each arm  $i \in \{1, \dots, n\}$ 
3:   Play arm  $i_t = \arg \max_i \hat{\mu}_{i,t}$ 
4:   Observe reward  $r_{i_t,t}$  and update posterior  $P_{i,t}$ 
5: end for

```

---

Recall the probability density function of the Beta distribution,  $\text{Beta}(\alpha, \beta)$ , for any  $x \in [0, 1]$  is

$$p(x) = \frac{(\alpha + \beta - 1)!}{(\alpha - 1)!(\beta - 1)!} x^{\alpha-1} (1-x)^{\beta-1}.$$

We also know, for any  $p, q \geq 1$ ,

$$\int_0^1 x^p (1-x)^q dx = \frac{(p+q+1)!}{(p-1)!(q-1)!}.$$

1. **[15 points]** Until time-step  $t$ , suppose arm  $i \in \{1, \dots, n\}$  is pulled  $N_{i,t}$  times and its total observed reward is

$$S_{i,t} := \sum_{u \leq t: i_u = i} r_{i,u},$$

where  $i_u$  represents the arm chosen at time-step  $u$ . Find  $P_{i,t}$ , the posterior distribution of  $\mu_i$ , given the Beta prior as described above and observations on the rewards until time-step  $t$ . (Hint: Compute the posterior for the first time-step. Use this recursively for the following time-steps.)

To find the posterior distribution  $P_{i,t}$  of the unknown mean reward  $\mu_i$  of arm  $i$  at time-step  $t$ , we begin with the prior and apply Bayesian updates.

Initially, the prior distribution of  $\mu_i$  is:

$$P_{i,0} = \text{Beta}(1, 1)$$

Suppose arm  $i$  has been pulled  $N_{i,t}$  times by time  $t$ , and the total observed reward is:

$$S_{i,t} := \sum_{u \leq t: i_u = i} r_{i,u}$$

Each reward  $r_{i,u} \sim \text{Bernoulli}(\mu_i)$ , and since the Beta distribution is a conjugate prior for the Bernoulli likelihood, the posterior distribution after observing  $S_{i,t}$  successes and  $N_{i,t} - S_{i,t}$  failures becomes:

$$P_{i,t} = \text{Beta}(1 + S_{i,t}, 1 + N_{i,t} - S_{i,t})$$

This result can be obtained recursively. After each new observation of reward  $r_{i,t}$  at time-step  $t$ , the posterior is updated as:

$$P_{i,t} = \text{Beta}(\alpha_{i,t}, \beta_{i,t})$$

where

$$\alpha_{i,t} = 1 + \sum_{u \leq t: i_u = i} r_{i,u} = 1 + S_{i,t}, \quad \beta_{i,t} = 1 + \sum_{u \leq t: i_u = i} (1 - r_{i,u}) = 1 + N_{i,t} - S_{i,t}$$

Hence, the posterior distribution  $P_{i,t}$  at time-step  $t$  is:

$P_{i,t} = \text{Beta}(1 + S_{i,t}, 1 + N_{i,t} - S_{i,t})$

2. **[5 points]** Compute the mean and variance of the posterior distribution of  $\mu_i$  found in part (a).

Let  $P_{i,t} = \text{Beta}(\alpha, \beta)$  where

$$\alpha = 1 + S_{i,t}, \quad \beta = 1 + N_{i,t} - S_{i,t}$$

Then the **mean** of  $\mu_i$  is:

$$\mathbb{E}[\mu_i] = \frac{1 + S_{i,t}}{2 + N_{i,t}}$$

And the **variance** of  $\mu_i$  is:

$$\text{Var}[\mu_i] = \frac{(1 + S_{i,t})(1 + N_{i,t} - S_{i,t})}{(2 + N_{i,t})^2(3 + N_{i,t})}$$

3. [5 points] Using the computations in part (b), explain how TS balances exploration and exploitation.

Thompson Sampling (TS) naturally balances exploration and exploitation by using random samples drawn from the posterior distribution of each arm's expected reward.

From part (b), we know that the posterior distribution of the mean reward  $\mu_i$  for arm  $i$  at time  $t$  is:

$$P_{i,t} = \text{Beta}(1 + S_{i,t}, 1 + N_{i,t} - S_{i,t})$$

with mean:

$$\mathbb{E}[\mu_i] = \frac{1 + S_{i,t}}{2 + N_{i,t}}$$

and variance:

$$\text{Var}[\mu_i] = \frac{(1 + S_{i,t})(1 + N_{i,t} - S_{i,t})}{(2 + N_{i,t})^2(3 + N_{i,t})}$$

**Exploitation:** When an arm is pulled frequently (i.e.,  $N_{i,t}$  is large), the variance of its posterior distribution becomes small. As a result, the samples  $\hat{\mu}_{i,t}$  drawn from this posterior concentrate near the true mean  $\mu_i$ , leading TS to consistently select arms with higher estimated reward.

**Exploration:** For arms that have been pulled fewer times (i.e., small  $N_{i,t}$ ), the variance of their posterior remains large. This increases the probability of drawing a high sample  $\hat{\mu}_{i,t}$  by chance, which may temporarily make that arm appear better than others. This leads TS to occasionally explore uncertain arms, even if their empirical average is lower.

## 4 Gridworld

Consider the following grid environment. Starting from any unshaded square, you can move up, down, left, or right. Actions are deterministic and always succeed (e.g., going left from state 16 goes to state 15) unless they will cause the agent to run into a wall. The thicker edges indicate walls, and attempting to move in the direction of a wall results in staying in the same square (e.g., going in any direction other than left from state 16 stays in 16).

Taking any action from the target square with cheese (no. 11) earns a reward of  $r_g$  (so  $r(11, a) = r_g \forall a$ ) and ends the episode. Taking any action from the square of the cat (no. 6) earns a reward of  $r_r$  (so  $r(6, a) = r_r \forall a$ ) and ends the episode. Otherwise, from every other square, taking any action is associated with a reward  $r_s \in \{-1, 0, +1\}$  (even if the action results in the agent staying in the same square).

Assume the discount factor  $\gamma = 1$ ,  $r_g = +10$ , and  $r_r = -1000$  unless otherwise specified.

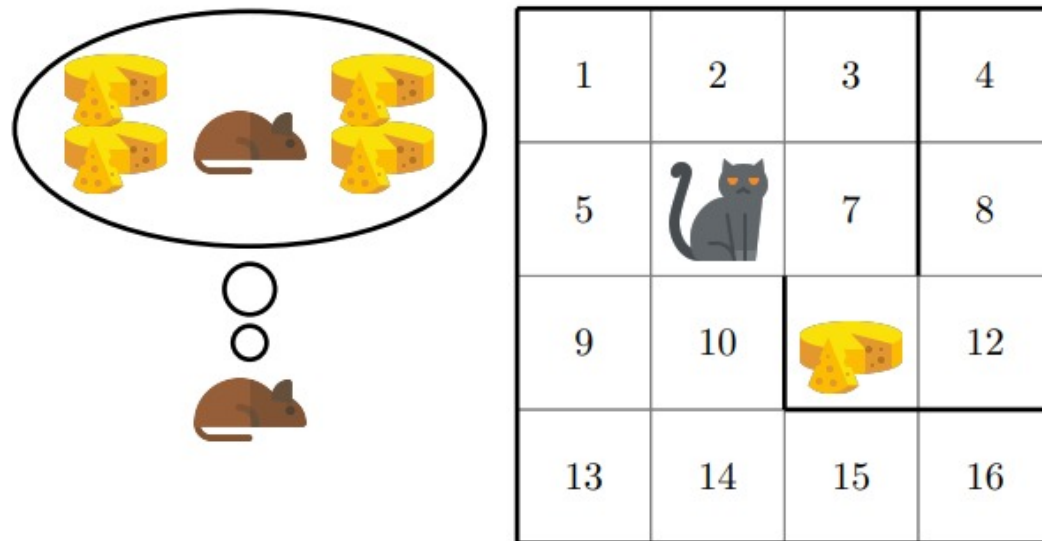


Figure 2: Gridworld Environment

- Let  $r_s = 0$ . In the policy below, the arrow in each box denotes the (deterministic) action that should be taken in each state. Evaluate the following policy and show the corresponding value (i.e., cumulative reward) for every state (square).

→	→	↓	↓
→	→	↓	↓
↑	↑	↓	←
↑	↑	←	←

- State 11 is the cheese (terminal, reward  $r_g = 10$ ).
  - State 6 is the cat (terminal, reward  $r_r = -1000$ ).
  - Non-terminal states have reward  $r_s = 0$ .
  - Discount factor  $\gamma = 1$ .
  - Actions are deterministic: up, down, left, right, unless blocked by walls (e.g., from state 16, only left to 15 is possible; other moves stay in 16).
  - The policy specifies actions as given.
- We need to compute the value function  $V^\pi(s)$ , the expected cumulative reward starting from state  $s$  under policy  $\pi$ . Since  $\gamma = 1$ ,  $r_s = 0$ , and actions are deterministic,  $V^\pi(s)$  is the terminal reward (10 or -1000) reached by following the policy.
- Policy Actions:**
  - Row 1:  $1 \rightarrow 2$ ,  $2 \rightarrow 3$ ,  $3 \downarrow 7$ ,  $4 \downarrow 8$
  - Row 2:  $5 \rightarrow 6$ ,  $6 \rightarrow 7$ ,  $7 \downarrow 11$ ,  $8 \downarrow 12$
  - Row 3:  $9 \uparrow 5$ ,  $10 \uparrow 6$ ,  $11 \downarrow 15$ ,  $12 \leftarrow 11$
  - Row 4:  $13 \uparrow 9$ ,  $14 \uparrow 10$ ,  $15 \leftarrow 14$ ,  $16 \leftarrow 15$
- Evaluate Paths:**
  - State 1:** Path:  $1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$ . Rewards:  $0 + 0 + 0 + 0 + 10 = 10$ . Thus,  $V^\pi(1) = 10$ .
  - State 2:** Path:  $2 \rightarrow 3 \downarrow 7 \downarrow 11$ . Rewards:  $0 + 0 + 0 + 10 = 10$ . Thus,  $V^\pi(2) = 10$ .
  - State 3:** Path:  $3 \downarrow 7 \downarrow 11$ . Rewards:  $0 + 0 + 10 = 10$ . Thus,  $V^\pi(3) = 10$ .
  - State 4:** Path:  $4 \downarrow 8 \downarrow 12 \leftarrow 11$ . Rewards:  $0 + 0 + 0 + 10 = 10$ . Thus,  $V^\pi(4) = 10$ .
  - State 5:** Path:  $5 \rightarrow 6$ . Rewards:  $0 - 1000 = -1000$ . Thus,  $V^\pi(5) = -1000$ .
  - State 6:** Terminal (cat). Reward:  $-1000$ . Thus,  $V^\pi(6) = -1000$ .
  - State 7:** Path:  $7 \downarrow 11$ . Rewards:  $0 + 10 = 10$ . Thus,  $V^\pi(7) = 10$ .
  - State 8:** Path:  $8 \downarrow 12 \leftarrow 11$ . Rewards:  $0 + 0 + 10 = 10$ . Thus,  $V^\pi(8) = 10$ .
  - State 9:** Path:  $9 \uparrow 5 \rightarrow 6$ . Rewards:  $0 + 0 - 1000 = -1000$ . Thus,  $V^\pi(9) = -1000$ .
  - State 10:** Path:  $10 \uparrow 6$ . Rewards:  $0 - 1000 = -1000$ . Thus,  $V^\pi(10) = -1000$ .
  - State 11:** Terminal (cheese). Reward:  $10$ . Thus,  $V^\pi(11) = 10$ .
  - State 12:** Path:  $12 \leftarrow 11$ . Rewards:  $0 + 10 = 10$ . Thus,  $V^\pi(12) = 10$ .
  - State 13:** Path:  $13 \uparrow 9 \uparrow 5 \rightarrow 6$ . Rewards:  $0 + 0 + 0 - 1000 = -1000$ . Thus,  $V^\pi(13) = -1000$ .
  - State 14:** Path:  $14 \uparrow 10 \uparrow 6$ . Rewards:  $0 + 0 - 1000 = -1000$ . Thus,  $V^\pi(14) = -1000$ .

- **State 15:** Path:  $15 \leftarrow 14 \uparrow 10 \uparrow 6$ . Rewards:  $0 + 0 + 0 - 1000 = -1000$ . Thus,  $V^\pi(15) = -1000$ .
- **State 16:** Path:  $16 \leftarrow 15 \leftarrow 14 \uparrow 10 \uparrow 6$ . Rewards:  $0 + 0 + 0 + 0 - 1000 = -1000$ . Thus,  $V^\pi(16) = -1000$ .

- **Values:**

10	10	10	10
-1000	-1000	10	10
-1000	-1000	10	10
-1000	-1000	-1000	-1000

The values under the policy are:

- States 1, 2, 3, 4, 7, 8, 11, 12:  $V^\pi = 10$
- States 5, 6, 9, 10, 13, 14, 15, 16:  $V^\pi = -1000$

2. Define the value of  $r_s$  that would cause the optimal policy to return the shortest path to the cheese square (no. 11); note that your value for  $r_s$  should lie in  $\{-1, 0, +1\}$ . Using this  $r_s$ , find the optimal value for each square and the optimal policy.

- **Setup:**

- State 11: Cheese (terminal, reward  $r_g = 10$ ).
- State 6: Cat (terminal, reward  $r_r = -1000$ ).
- Non-terminal states: Reward  $r_s \in \{-1, 0, +1\}$ .
- Discount factor  $\gamma = 1$ .
- Actions: Up, down, left, right; deterministic unless blocked by walls (e.g., from state 16, only left to 15 is possible).
- Walls: State 16 allows only left to 15; standard boundaries.

- **Choosing  $r_s$ :** To ensure the optimal policy takes the shortest path to state 11:

- $r_s = -1$ : Penalizes each step, encouraging the shortest path.
- $r_s = 0$ : No penalty, so all paths to state 11 yield 10, not necessarily the shortest.
- $r_s = +1$ : Rewards steps, encouraging longer paths.

Thus,  $r_s = -1$ .

- **Optimal Value Function:** With  $r_s = -1$ ,  $\gamma = 1$ , for non-terminal states:

$$V^*(s) = \max_a [-1 + V^*(s')]$$

Terminal:  $V^*(11) = 10$ ,  $V^*(6) = -1000$ . Values are  $V^*(s) = 10 - k$ , where  $k$  is the number of steps in the shortest path to state 11.

- **Shortest Paths:**

- State 1:  $1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (4 steps).  $V^*(1) = 10 - 4 = 6$ . Action:  $\rightarrow$ .
- State 2:  $2 \rightarrow 3 \downarrow 7 \downarrow 11$  (3 steps).  $V^*(2) = 10 - 3 = 7$ . Action:  $\rightarrow$ .
- State 3:  $3 \downarrow 7 \downarrow 11$  (2 steps).  $V^*(3) = 10 - 2 = 8$ . Action:  $\downarrow$ .
- State 4:  $4 \downarrow 8 \downarrow 12 \leftarrow 11$  (3 steps).  $V^*(4) = 10 - 3 = 7$ . Action:  $\downarrow$ .
- State 5:  $5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (5 steps).  $V^*(5) = 10 - 5 = 5$ . Action:  $\uparrow$ .
- State 6: Terminal.  $V^*(6) = -1000$ .
- State 7:  $7 \downarrow 11$  (1 step).  $V^*(7) = 10 - 1 = 9$ . Action:  $\downarrow$ .
- State 8:  $8 \downarrow 12 \leftarrow 11$  (2 steps).  $V^*(8) = 10 - 2 = 8$ . Action:  $\downarrow$ .
- State 9:  $9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (6 steps).  $V^*(9) = 10 - 6 = 4$ . Action:  $\uparrow$ .
- State 10:  $10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (7 steps).  $V^*(10) = 10 - 7 = 3$ . Action:  $\leftarrow$ .
- State 11: Terminal.  $V^*(11) = 10$ .
- State 12:  $12 \leftarrow 11$  (1 step).  $V^*(12) = 10 - 1 = 9$ . Action:  $\leftarrow$ .
- State 13:  $13 \uparrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (7 steps).  $V^*(13) = 10 - 7 = 3$ . Action:  $\uparrow$ .
- State 14:  $14 \uparrow 10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (8 steps).  $V^*(14) = 10 - 8 = 2$ . Action:  $\uparrow$ .
- State 15:  $15 \leftarrow 14 \uparrow 10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (9 steps).  $V^*(15) = 10 - 9 = 1$ . Action:  $\leftarrow$ .
- State 16:  $16 \leftarrow 15 \leftarrow 14 \uparrow 10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (10 steps).  $V^*(16) = 10 - 10 = 0$ . Action:  $\leftarrow$ .

- **Optimal Policy:**

$\rightarrow$	$\rightarrow$	$\downarrow$	$\downarrow$
$\uparrow$	T	$\downarrow$	$\downarrow$
$\uparrow$	$\leftarrow$	T	$\leftarrow$
$\uparrow$	$\uparrow$	$\leftarrow$	$\leftarrow$

where T denotes terminal states.

- **Values:**

6	7	8	7
5	-1000	9	8
4	3	10	9
3	2	1	0

**Answer:**

- $r_s = -1$ .
- Optimal values:

- State 1: 6, State 2: 7, State 3: 8, State 4: 7, State 5: 5, State 6: -1000,
- State 7: 9, State 8: 8, State 9: 4, State 10: 3, State 11: 10, State 12: 9,
- State 13: 3, State 14: 2, State 15: 1, State 16: 0

- Optimal policy: As shown in the table above.

3. Let's refer to the value function derived in (b) as  $V_{\text{old}}^{\pi_g}$  and the policy as  $\pi_g$ . Suppose we are now in a new gridworld where all the rewards ( $r_s$ ,  $r_g$ , and  $r_r$ ) have +2 added to them. Consider still following the policy  $\pi_g$ , which is optimal for the original gridworld. What will the new values  $V_{\text{new}}^{\pi_g}$  be in this second gridworld?

- **Setup:**

- Original rewards:  $r_s = -1$ ,  $r_g = 10$  (state 11),  $r_r = -1000$  (state 6).
- New rewards:  $r_s = -1 + 2 = 1$ ,  $r_g = 10 + 2 = 12$ ,  $r_r = -1000 + 2 = -998$ .
- Discount factor  $\gamma = 1$ .
- Policy  $\pi_g$ :

→	→	↓	↓
↑	T	↓	↓
↑	←	T	←
↑	↑	←	←

- **Objective:** Compute  $V_{\text{new}}(s)$  under  $\pi_g$  with new rewards.

- **Compute  $V_{\text{new}}(s)$ :** Follow  $\pi_g$  to state 11, summing rewards ( $r_s = 1$ ,  $r_g = 12$ ).

- State 1:  $1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (4 steps). Rewards:  $1 + 1 + 1 + 1 + 12 = 16$ .  $V_{\text{new}}(1) = 16$ .
- State 2:  $2 \rightarrow 3 \downarrow 7 \downarrow 11$  (3 steps). Rewards:  $1 + 1 + 1 + 12 = 15$ .  $V_{\text{new}}(2) = 15$ .
- State 3:  $3 \downarrow 7 \downarrow 11$  (2 steps). Rewards:  $1 + 1 + 12 = 14$ .  $V_{\text{new}}(3) = 14$ .
- State 4:  $4 \downarrow 8 \downarrow 12 \leftarrow 11$  (3 steps). Rewards:  $1 + 1 + 1 + 12 = 15$ .  $V_{\text{new}}(4) = 15$ .
- State 5:  $5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (5 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 12 = 17$ .  $V_{\text{new}}(5) = 17$ .
- State 6: Terminal.  $V_{\text{new}}(6) = -998$ .
- State 7:  $7 \downarrow 11$  (1 step). Rewards:  $1 + 12 = 13$ .  $V_{\text{new}}(7) = 13$ .
- State 8:  $8 \downarrow 12 \leftarrow 11$  (2 steps). Rewards:  $1 + 1 + 12 = 14$ .  $V_{\text{new}}(8) = 14$ .
- State 9:  $9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (6 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 1 + 12 = 18$ .  $V_{\text{new}}(9) = 18$ .
- State 10:  $10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (7 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 1 + 1 + 12 = 19$ .  $V_{\text{new}}(10) = 19$ .
- State 11: Terminal.  $V_{\text{new}}(11) = 12$ .
- State 12:  $12 \leftarrow 11$  (1 step). Rewards:  $1 + 12 = 13$ .  $V_{\text{new}}(12) = 13$ .
- State 13:  $13 \uparrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (7 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 1 + 1 + 12 = 19$ .  $V_{\text{new}}(13) = 19$ .
- State 14:  $14 \uparrow 10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (8 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 12 = 20$ .  $V_{\text{new}}(14) = 20$ .
- State 15:  $15 \leftarrow 14 \uparrow 10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (9 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 12 = 21$ .  $V_{\text{new}}(15) = 21$ .
- State 16:  $16 \leftarrow 15 \leftarrow 14 \uparrow 10 \leftarrow 9 \uparrow 5 \uparrow 1 \rightarrow 2 \rightarrow 3 \downarrow 7 \downarrow 11$  (10 steps). Rewards:  $1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 12 = 22$ .  $V_{\text{new}}(16) = 22$ .

- **Values:**

16	15	14	15
17	-998	13	14
18	19	12	13
19	20	21	22

**Answer:** The new values under policy  $\pi_g$  are:

- State 1: 16, State 2: 15, State 3: 14, State 4: 15, State 5: 17, State 6: -998,
- State 7: 13, State 8: 14, State 9: 18, State 10: 19, State 11: 12, State 12: 13,
- State 13: 19, State 14: 20, State 15: 21, State 16: 22



## 5 Markov Decision Process as Linear Program

Consider a Markov Decision Process (MDP)  $M = \{S, A, P, r, \gamma\}$  where  $S$  is the state space,  $A$  is the action space,  $P$  is the transition kernel representing the probability of transition to state  $s'$  from state  $s$  when action  $a$  is taken as  $P(s'|s, a)$ ,  $r$  is the reward function, and  $\gamma$  is the discount factor. Consider the linear program:

$$\min_{V \in \mathbb{R}^{|S|}} \sum_{s \in S} \rho(s) V(s)$$

subject to

$$V(s) \geq r(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V(s') \quad \forall s \in S, a \in A,$$

where  $\rho(s) \in \mathbb{R}_+ \forall s \in S$  and denote its solution as  $V^* \in \mathbb{R}^{|S|}$ .

1. Show that the dual formulation of the above linear program can be written as

$$\max_{x \in \mathbb{R}^{|S| \times |A|}} \sum_{s \in S} \sum_{a \in A} r(s, a) x(s, a)$$

subject to

$$\begin{aligned} \sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S} \sum_{a' \in A} P(s|s', a') x(s', a') &= \rho(s) \quad \forall s \in S, \\ x(s, a) &\geq 0 \quad \forall s \in S, a \in A. \end{aligned}$$

To derive the dual, rewrite the constraint in standard form:

$$V(s) - \gamma \sum_{s' \in S} P(s' | s, a) V(s') \geq r(s, a).$$

Associate a dual variable  $x(s, a) \geq 0$  with each constraint. The Lagrangian is:

$$L(V, x) = \sum_{s \in S} \rho(s) V(s) + \sum_{s \in S} \sum_{a \in A} x(s, a) \left[ r(s, a) - V(s) + \gamma \sum_{s' \in S} P(s' | s, a) V(s') \right].$$

Collect terms involving  $V(s)$ :

$$L(V, x) = \sum_{s \in S} V(s) \left[ \rho(s) - \sum_{a \in A} x(s, a) + \gamma \sum_{s' \in S} \sum_{a' \in A} P(s | s', a') x(s', a') \right] + \sum_{s \in S} \sum_{a \in A} x(s, a) r(s, a).$$

For the dual, set the coefficient of  $V(s)$  to zero:

$$\rho(s) - \sum_{a \in A} x(s, a) + \gamma \sum_{s' \in S} \sum_{a' \in A} P(s | s', a') x(s', a') = 0 \quad \forall s \in S.$$

Thus, the dual constraint is:

$$\sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S, a' \in A} P(s | s', a') x(s', a') = \rho(s) \quad \forall s \in S.$$

The dual objective is to maximize the constant term:

$$\max \sum_{s \in S} \sum_{a \in A} x(s, a) r(s, a).$$

Therefore, the dual linear program is:

$$\begin{aligned} &\max_{x \in \mathbb{R}^{|S| \times |A|}} \sum_{s \in S, a \in A} r(s, a) x(s, a) \\ \text{s.t.} \quad &\sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S, a' \in A} P(s | s', a') x(s', a') = \rho(s) \quad \forall s \in S, \\ &x(s, a) \geq 0 \quad \forall s \in S, a \in A. \end{aligned}$$

2. Denote the optimal solution to the dual problem as  $x^* \in \mathbb{R}^{|S| \times |A|}$ . Due to strong duality and complementary slackness, we have

$$x^*(s, a) \left( V^*(s) - \gamma \sum_{s' \in S} P(s'|s, a) V^*(s') - r(s, a) \right) = 0 \quad \forall s \in S, a \in A.$$

Now, show that the optimal policy  $\pi^*(\cdot|s)$ ,  $\forall s \in S$ , can be derived as

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in A} x^*(s, a), \\ 0 & \text{else.} \end{cases}$$

Given the optimal dual solution  $x^*(s, a)$  and primal solution  $V^*(s)$ , complementary slackness gives:

$$x^*(s, a) \left( V^*(s) - \gamma \sum_{s' \in S} P(s' | s, a) V^*(s') - r(s, a) \right) = 0 \quad \forall s \in S, a \in A.$$

If  $x^*(s, a) > 0$ , the constraint is tight:

$$V^*(s) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a) V^*(s').$$

Since  $V^*(s)$  is the optimal value function, it satisfies:

$$V^*(s) = \max_{a' \in \mathcal{A}} \left\{ r(s, a') + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a') V^*(s') \right\}.$$

Thus,  $x^*(s, a) > 0$  implies action  $a$  achieves the maximum, so it is optimal. Define the policy:

$$\pi^*(a \mid s) = \begin{cases} 1 & \text{if } a = \arg \max_{a' \in \mathcal{A}} x^*(s, a'), \\ 0 & \text{else.} \end{cases}$$

If  $a^* = \arg \max_{a'} x^*(s, a')$  and  $x^*(s, a^*) > 0$ , then  $a^*$  is optimal, satisfying the Bellman equation. If  $x^*(s, a) = 0$ , action  $a$  is suboptimal or not selected. Since an optimal policy can be deterministic, choosing  $a^* = \arg \max x^*(s, a)$  ensures optimality.

Thus, the optimal policy is:

$$\pi^*(a \mid s) = \begin{cases} 1 & \text{if } a = \arg \max_{a' \in \mathcal{A}} x^*(s, a'), \\ 0 & \text{else.} \end{cases}$$