# Lab3

## Oscar Hoffmann

## 2024-10-05

### 2.1 Q-Learning

Complete the Q-learning algorithm for a grid-world environment where the agent navigates using four actions (up, down, left, right) and faces a probability of slipping. The agent must learn the optimal policy using a model-free approach, as the transition model is unknown.

```r
rm(list =ls())
set.seed(1234)
library(ggplot2)

arrows <- c("^", ">", "v", "<")
action_deltas <- list(c(1,0), # up
                      c(0,1), # right
                      c(-1,0), # down
                      c(0,-1)) # left

vis_environment <- function(iterations=0, epsilon = 0.5, alpha = 0.1, gamma = 0.95, beta = 0){

  df <- expand.grid(x=1:H,y=1:W)
  foo <- mapply(function(x,y) ifelse(reward_map[x,y] == 0,q_table[x,y,1],NA),df$x,df$y)
  df$val1 <- as.vector(round(foo, 2))
  foo <- mapply(function(x,y) ifelse(reward_map[x,y] == 0,q_table[x,y,2],NA),df$x,df$y)
  df$val2 <- as.vector(round(foo, 2))
  foo <- mapply(function(x,y) ifelse(reward_map[x,y] == 0,q_table[x,y,3],NA),df$x,df$y)
  df$val3 <- as.vector(round(foo, 2))
  foo <- mapply(function(x,y) ifelse(reward_map[x,y] == 0,q_table[x,y,4],NA),df$x,df$y)
  df$val4 <- as.vector(round(foo, 2))
  foo <- mapply(function(x,y)
    ifelse(reward_map[x,y] == 0,arrows[GreedyPolicy(x,y)],reward_map[x,y]),df$x,df$y)
  df$val5 <- as.vector(foo)
  foo <- mapply(function(x,y) ifelse(reward_map[x,y] == 0,max(q_table[x,y,]),
                                     ifelse(reward_map[x,y]<0,NA,reward_map[x,y])),df$x,df$y)
  df$val6 <- as.vector(foo)

  print(ggplot(df,aes(x = y,y = x)) +
          scale_fill_gradient(low = "white", high = "green", na.value = "red", name = "") +
          geom_tile(aes(fill=val6)) +
          geom_text(aes(label = val1),size = 4,nudge_y = .35,na.rm = TRUE) +
          geom_text(aes(label = val2),size = 4,nudge_x = .35,na.rm = TRUE) +
          geom_text(aes(label = val3),size = 4,nudge_y = -.35,na.rm = TRUE) +
          geom_text(aes(label = val4),size = 4,nudge_x = -.35,na.rm = TRUE) +
          geom_text(aes(label = val5),size = 10) +
```

```r
            geom_tile(fill = 'transparent', colour = 'black') +
            ggtitle(paste("Q-table after ",iterations," iterations\n",
                          "(epsilon = ",epsilon,", alpha = ",alpha,"gamma = ",gamma,", beta = ",beta,")"))
            theme(plot.title = element_text(hjust = 0.5)) +
            scale_x_continuous(breaks = c(1:W),labels = c(1:W)) +
            scale_y_continuous(breaks = c(1:H),labels = c(1:H)))
}

GreedyPolicy <- function(x, y){
  # Get the Q-values for all actions at state (x, y)
  q_values <- q_table[x, y, ]

  # Find the max Q-value
  max_q <- max(q_values)

  # Identify all actions with maximum Q-value
  max_actions <- which(q_values == max_q)

  # Check and resolve ties
  if (length(max_actions) > 1) {
    action <- sample(max_actions, 1)
  } else {
    action <- max_actions
  }
  return(action)
}

EpsilonGreedyPolicy <- function(x, y, epsilon){
  # Generate a random numb
  rand_num <- runif(1)
  if (rand_num < epsilon){
   #select a random action
    action <- sample(1:4, 1)
  } else {
  # use the greedy policy
    action <- GreedyPolicy(x, y)
  }
  return(action)
}

transition_model <- function(x, y, action, beta){
  delta <- sample(-1:1, size = 1, prob = c(0.5*beta,1-beta,0.5*beta))
  final_action <- ((action + delta + 3) %% 4) + 1
  foo <- c(x,y) + unlist(action_deltas[final_action])
  foo <- pmax(c(1,1),pmin(foo,c(H,W)))
  return (foo)
}

q_learning <- function(start_state, epsilon = 0.5, alpha = 0.1, gamma = 0.95, beta = 0){

  # Perform one episode of Q-learning. The agent should move around in the
  # environment using the given transition model and update the Q-table.
  # The episode ends when the agent reaches a terminal state.
```

```r
#
# Args:
#   start_state: array with two entries, describing the starting position of the agent.
#   epsilon (optional): probability of acting randomly.
#   alpha (optional): learning rate.
#   gamma (optional): discount factor.
#   beta (optional): slipping factor.
#   reward_map (global variable): a HxW array containing the reward given at each state.
#   q_table (global variable): a HxWx4 array containing Q-values for each state-action pair.
#
# Returns:
#   reward: reward received in the episode.
#   correction: sum of the temporal difference correction terms over the episode.
#   q_table (global variable): Recall that R passes arguments by value. So, q_table being
#   a global variable can be modified with the superassigment operator <<-.

# Initialize state
x <- start_state[1]
y <- start_state[2]

# Initialize variables to track episode reward and TD corrections
episode_reward <- 0
episode_correction <- 0

repeat{
  # Choose an action A with epsilon-greedy policy
  action <- EpsilonGreedyPolicy(x, y, epsilon)

  # Observe next state S' & reward R after taking action A
  next_state <- transition_model(x, y, action, beta)
  x_new <- next_state[1]
  y_new <- next_state[2]
  R <- reward_map[x_new, y_new]

  # Get current Q-value Q(S, A)
  Q_SA <- q_table[x, y, action]

  # if R != 0 it means next state is terminal, this check if for when the terminal state
  # is the first action and the end of the loop hasn't been reached
  if (R != 0){
    max_QSAprime <- 0
  } else {
    # Compute max Q(S', a') over all possible actions a'
    max_QSAprime <- max(q_table[x_new, y_new, ])
  }

  # Calc TD correction term
  TD_correction <- R + gamma * max_QSAprime - Q_SA
  episode_correction <- episode_correction + TD_correction

  # Update Q(S, A)
  q_table[x, y, action] <<- Q_SA + alpha * TD_correction
```

```
    # Ep reward accumulated
    episode_reward <- episode_reward + R

    # Next state
    x <- x_new
    y <- y_new

    # Check if the episode has ended (terminal state)
    if (R != 0){
      return (c(episode_reward, episode_correction))
    }
  }
}
```

## 2.2 Environment A

Implement and visualize Q-learning in a 5×7 grid-world environment with specific rewards, and complete tasks including implementing greedy and epsilon-greedy policies, and running 10,000 episodes of Q-learning to analyze the learned policy and Q-table.
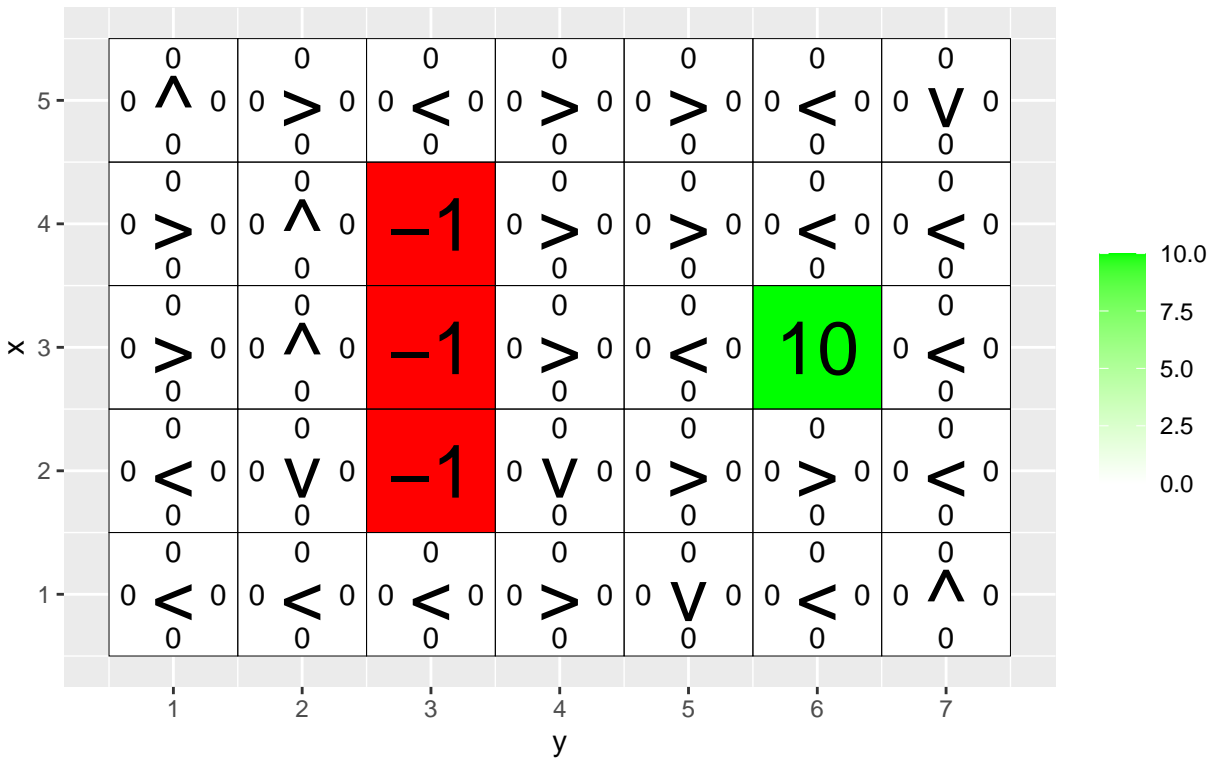
```
H <- 5
W <- 7

reward_map <- matrix(0, nrow = H, ncol = W)
reward_map[3,6] <- 10
reward_map[2:4,3] <- -1

q_table <- array(0,dim = c(H,W,4))

vis_environment()
```
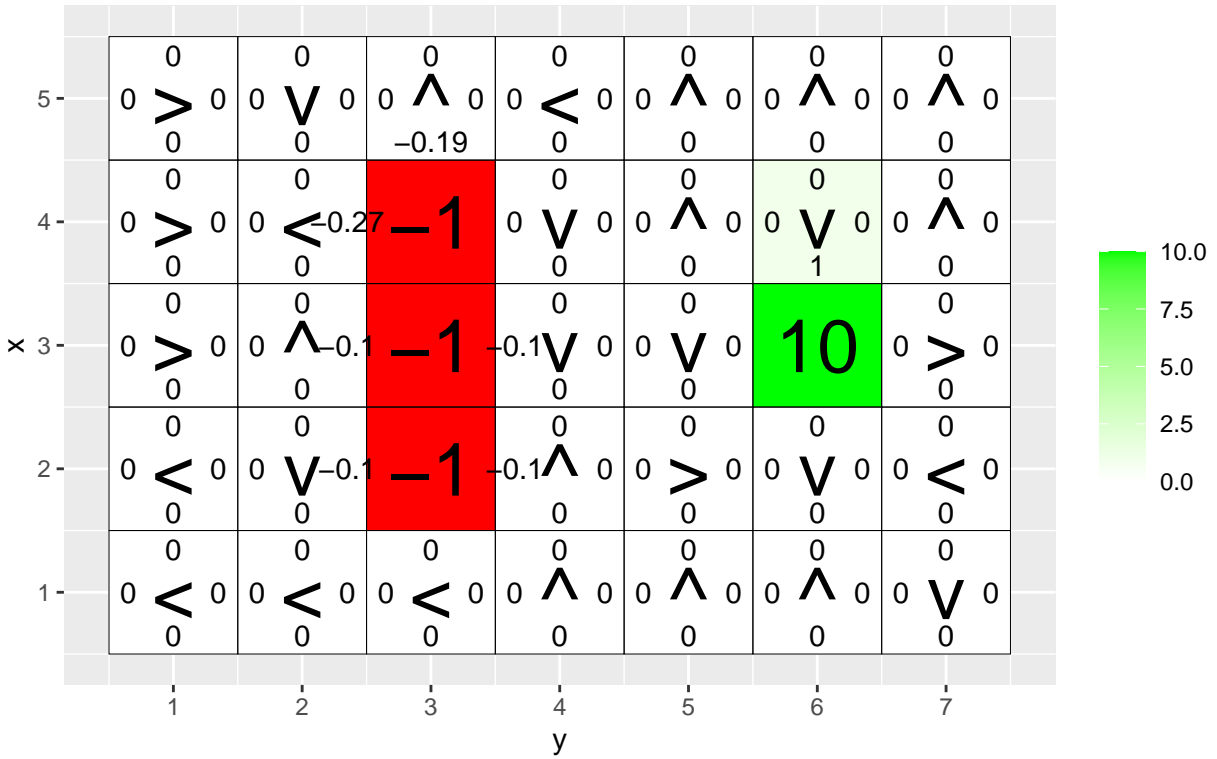
Q−table after 0 iterations
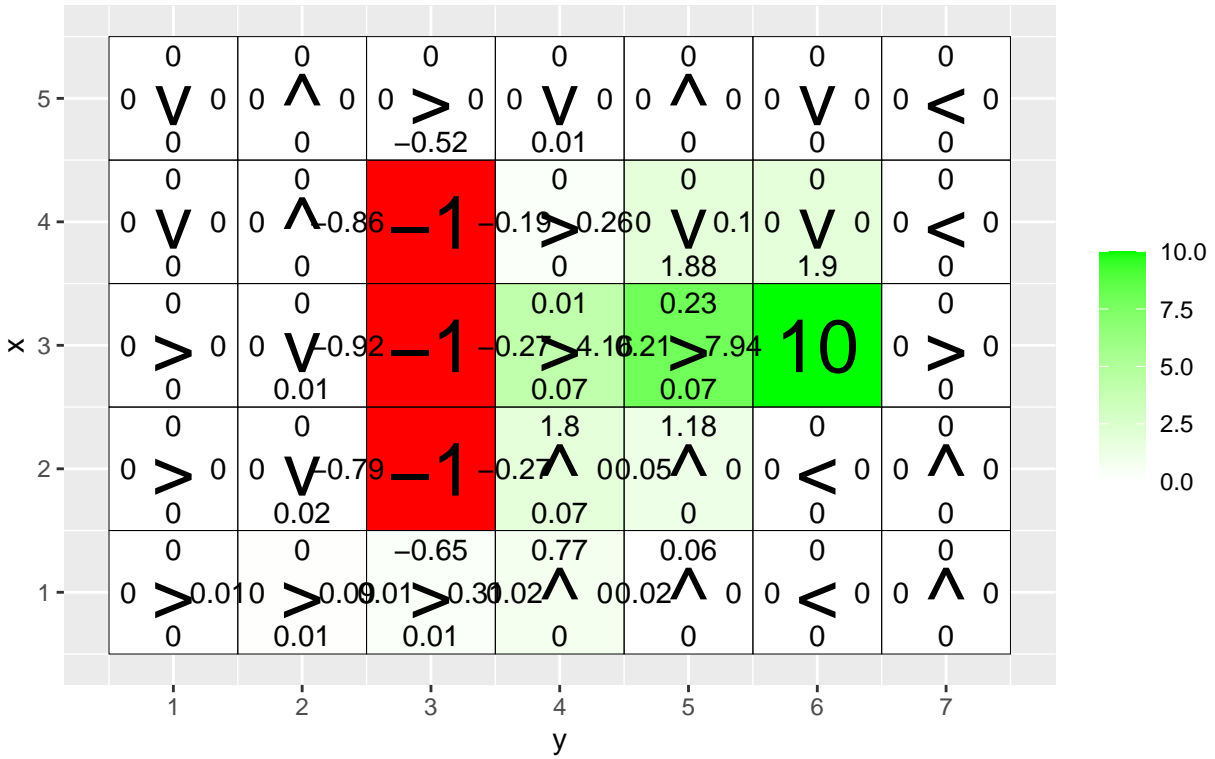(epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )

```
for(i in 1:10000){
  foo <- q_learning(start_state = c(3,1))

  if(any(i==c(10,100,1000,10000)))
    vis_environment(i)
}
```
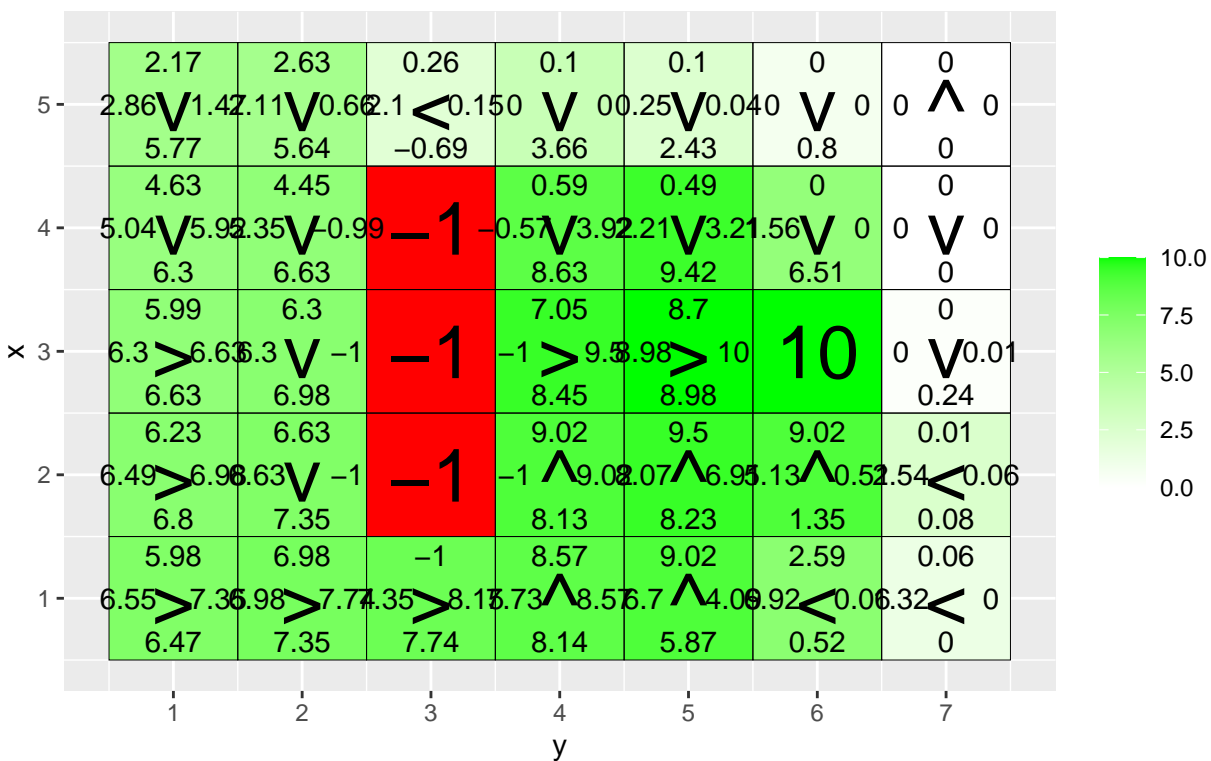
Q−table after 10 iterations
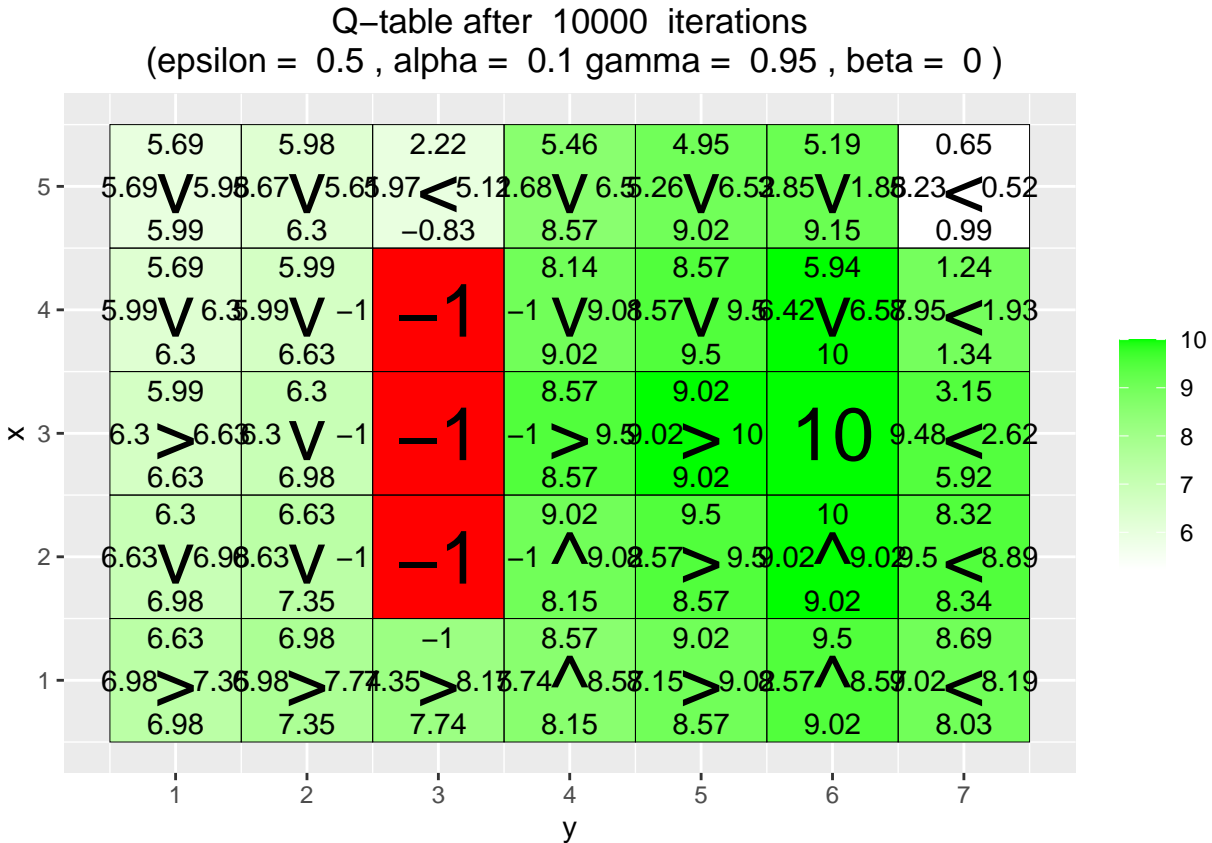(epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 5 | 0 > 0 / 0 | 0 V 0 / 0 | 0 ∧ 0 / −0.19 | 0 < 0 / 0 | 0 ∧ 0 / 0 | 0 ∧ 0 / 0 | 0 ∧ 0 / 0 |
| 4 | 0 > 0 / 0 | 0 < −0.27 / 0 | −1 | 0 V 0 / 0 | 0 ∧ 0 / 0 | 0 V 0 / 1 | 0 ∧ 0 / 0 |
| 3 | 0 > 0 / 0 | 0 ∧ −0.1 / 0 | −1 | −0.1 V 0 / 0 | 0 V 0 / 0 | 10 | 0 > 0 / 0 |
| 2 | 0 < 0 / 0 | 0 V −0.1 / 0 | −1 | −0.1 ∧ 0 / 0 | 0 > 0 / 0 | 0 V 0 / 0 | 0 < 0 / 0 |
| 1 | 0 < 0 / 0 | 0 < 0 / 0 | 0 < 0 / 0 | 0 ∧ 0 / 0 | 0 ∧ 0 / 0 | 0 ∧ 0 / 0 | 0 V 0 / 0 |

x

y

10.0
7.5
5.0
2.5
0.0

# Q−table after 100 iterations
## (epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )



7

Q–table after  1000  iterations
(epsilon =  0.5 , alpha =  0.1 gamma =  0.95 , beta =  0 )

x

5 -
| | 2.17 | 2.63 | 0.26 | 0.1 | 0.1 | 0 | 0 |
| 2.86 V 1.42 | 2.11 V 0.66 | 2.1 < 0.15 | 0 V 0 | 0.25 V 0.04 | 0 V 0 | 0 ^ 0 |
| 5.77 | 5.64 | −0.69 | 3.66 | 2.43 | 0.8 | 0 |

4 -
| 4.63 | 4.45 | | 0.59 | 0.49 | 0 | 0 |
| 5.04 V 5.92 | 2.35 V −0.99 | −1 | −0.57 V 3.9 | 2.21 V 3.2 | 1.56 V 0 | 0 V 0 |
| 6.3 | 6.63 | | 8.63 | 9.42 | 6.51 | 0 |

3 -
| 5.99 | 6.3 | | 7.05 | 8.7 | | 0 |
| 6.3 > 6.63 | 6.3 V −1 | −1 | −1 > 9.5 | 8.98 > 10 | 10 | 0 V 0.01 |
| 6.63 | 6.98 | | 8.45 | 8.98 | | 0.24 |

2 -
| 6.23 | 6.63 | | 9.02 | 9.5 | 9.02 | 0.01 |
| 6.49 > 6.98 | 6.63 V −1 | −1 | −1 ^ 9.02 | 2.07 ^ 6.9 | 5.13 ^ 0.5 | 2.54 < 0.06 |
| 6.8 | 7.35 | | 8.13 | 8.23 | 1.35 | 0.08 |

1 -
| 5.98 | 6.98 | −1 | 8.57 | 9.02 | 2.59 | 0.06 |
| 6.55 > 7.35 | 6.98 > 7.74 | 7.35 > 8.13 | 3.73 ^ 8.5 | 6.7 ^ 4.0 | 9.92 < 0.06 | 6.32 < 0 |
| 6.47 | 7.35 | 7.74 | 8.14 | 5.87 | 0.52 | 0 |

        1        2        3        4        5        6        7
                              y

10.0
7.5
5.0
2.5
0.0

## Q–table after 10000 iterations
### (epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )

**What has the agent learned after the first 10 episodes?**

The Q-table values are mostly zeros, with slight updates in some states around -1. The greedy policy (arrows on the grid) point randomly due to insufficient learning.

**Is the final greedy policy (after 10000 episodes) optimal for all states, i.e. not only for the initial state ? Why / Why not ?**

The agent has clearly learned to navigate towards the goal at position, as seen by the arrows pointing rightwards and downwards towards this cell. The policy is optimal for the most states that the agent has visited often during the episodes, however some areas are under explored leading to sub-optimal policies in those areas. An example is the states above the -1 penalty wall where it would be better to go above and not down again.

**Do the learned values in the Q-table reflect the fact that there are multiple paths (above and below the negative rewards) to get to the positive reward ? If not, what could be done to make it happen ?**

Mostly no, the learned values in the Q-table do not reflect the fact that there are multiple path that lead to the positive reward. The agent clearly shows a strong preference for the lower path over the upper path. This is likely due to the agent's early discovery of the upper path and subsequent exploitation of it.

To counteract this one could:

- increase exploration

- Reducing the learning rate
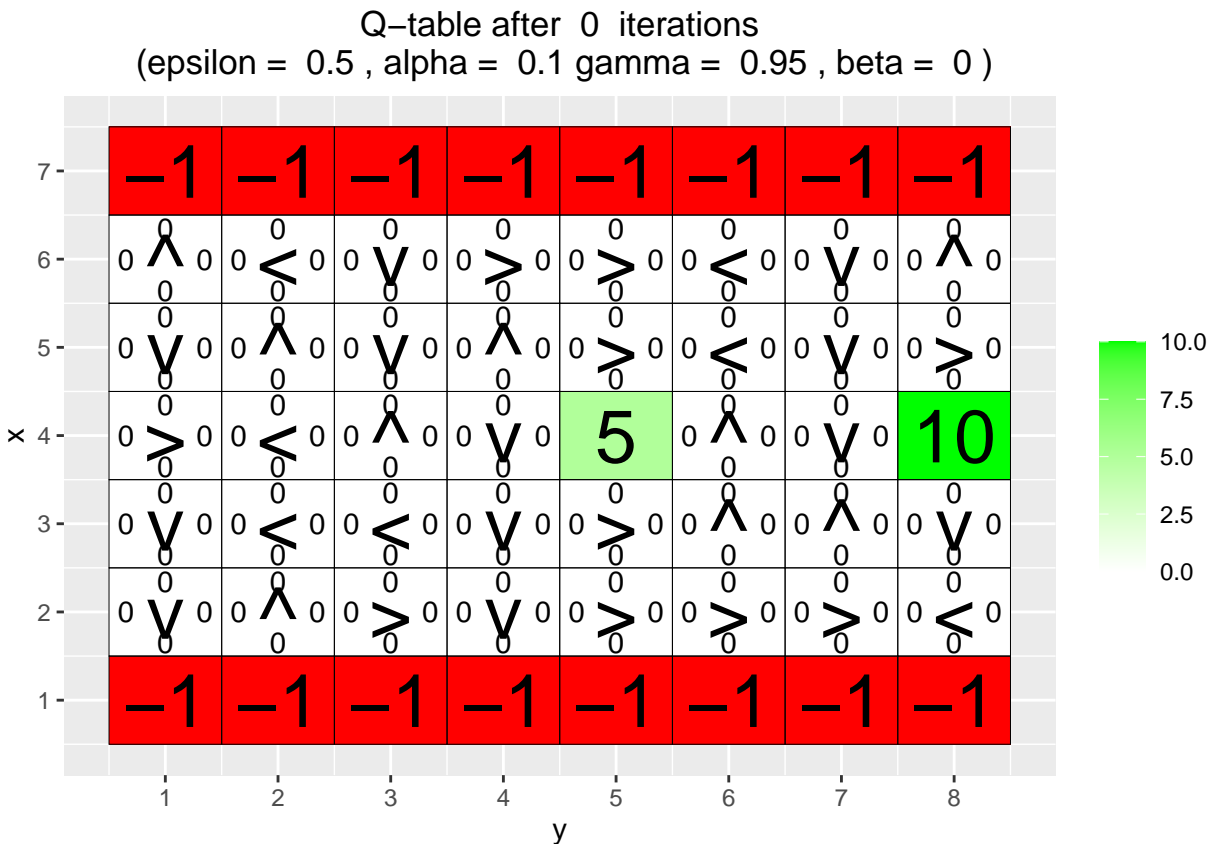
## 2.3 Environment B

Investigate how the parameters epsilon and gamma affect the learned policy in a 7×8 environment with negative and positive rewards by running 30,000 episodes of Q-learning with different parameter settings.

```r
H <- 7
W <- 8

reward_map <- matrix(0, nrow = H, ncol = W)
reward_map[1,] <- -1
reward_map[7,] <- -1
reward_map[4,5] <- 5
reward_map[4,8] <- 10

q_table <- array(0,dim = c(H,W,4))

vis_environment()
```



Q−table after 0 iterations
(epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )

```r
MovingAverage <- function(x, n){

  cx <- c(0,cumsum(x))
  rsum <- (cx[(n+1):length(cx)] - cx[1:(length(cx) - n)]) / n

  return (rsum)
}
```

```r
for(j in c(0.5,0.75,0.95)){
  q_table <- array(0,dim = c(H,W,4))
  reward <- NULL
  correction <- NULL

  for(i in 1:30000){
    # This is for epsilon = 0.5 standard
    foo <- q_learning(gamma = j, start_state = c(4,1))
    reward <- c(reward,foo[1])
    correction <- c(correction,foo[2])
  }

  vis_environment(i, gamma = j)
  plot(MovingAverage(reward,100),type = "l")
  plot(MovingAverage(correction,100),type = "l")
}
```
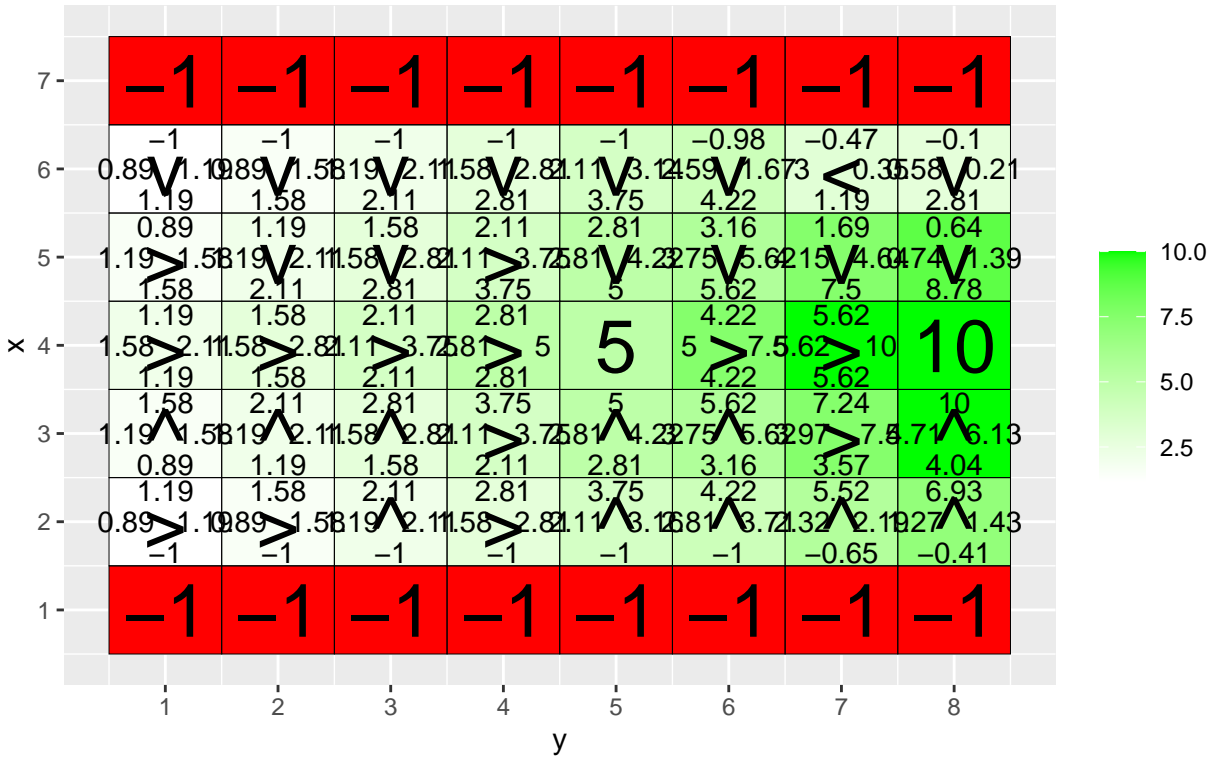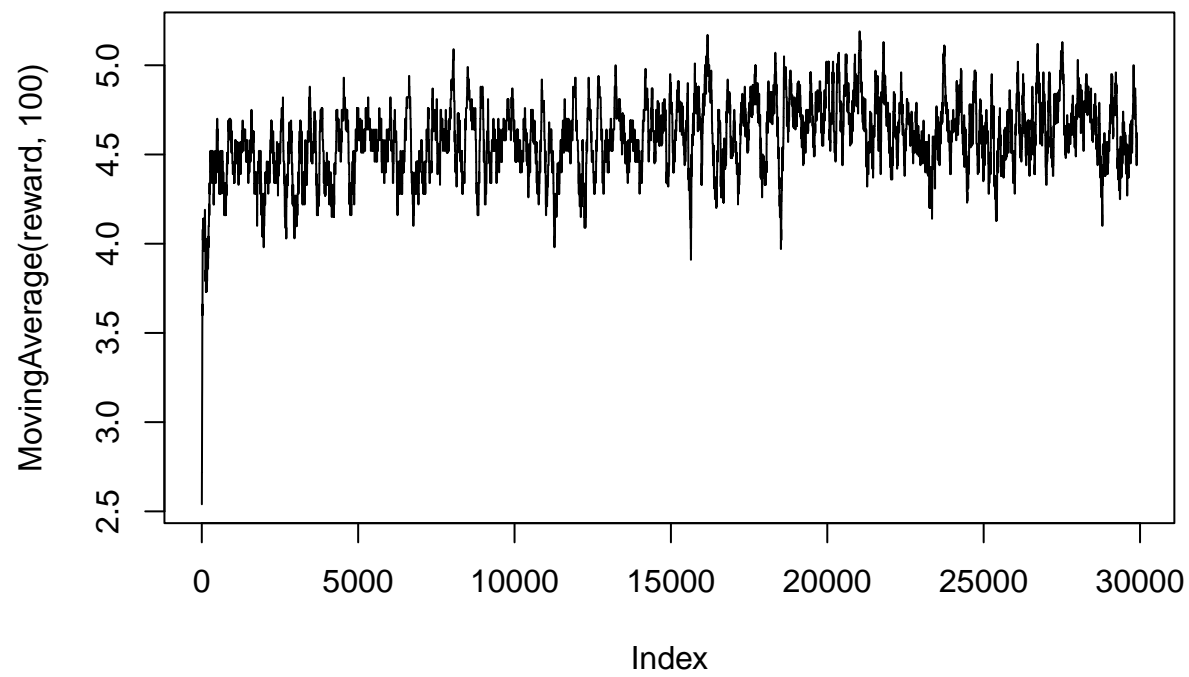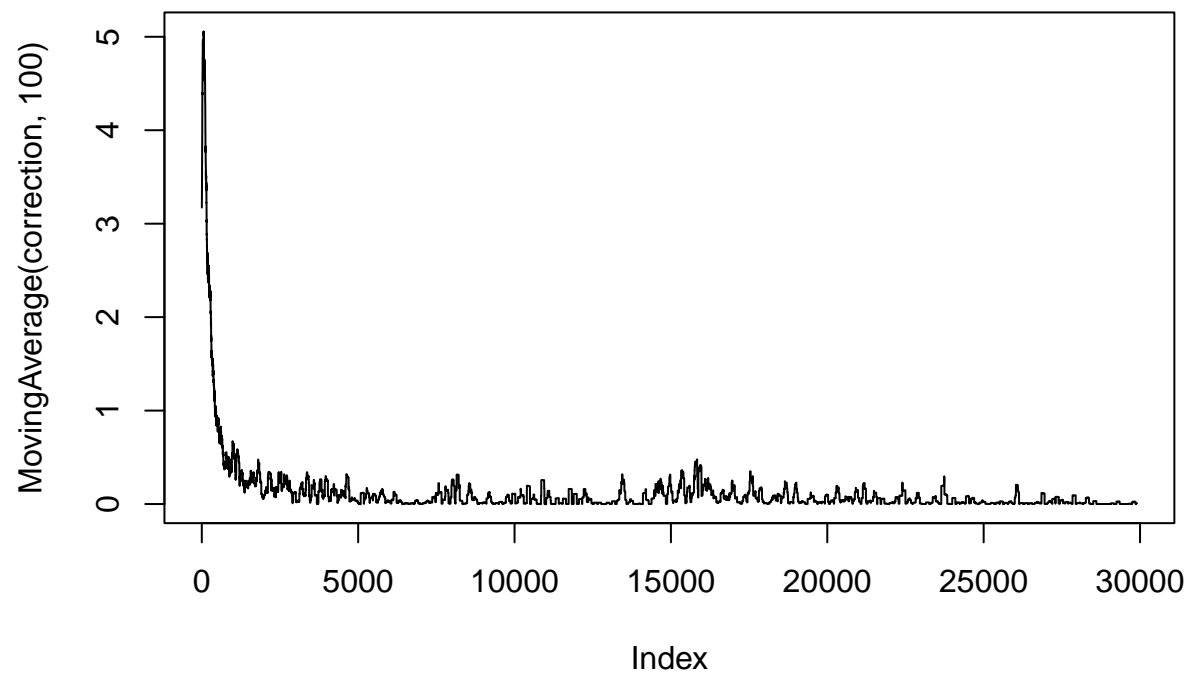


Q−table after 30000 iterations
(epsilon = 0.5 , alpha = 0.1 gamma = 0.5 , beta = 0 )

Q–table after 30000 iterations
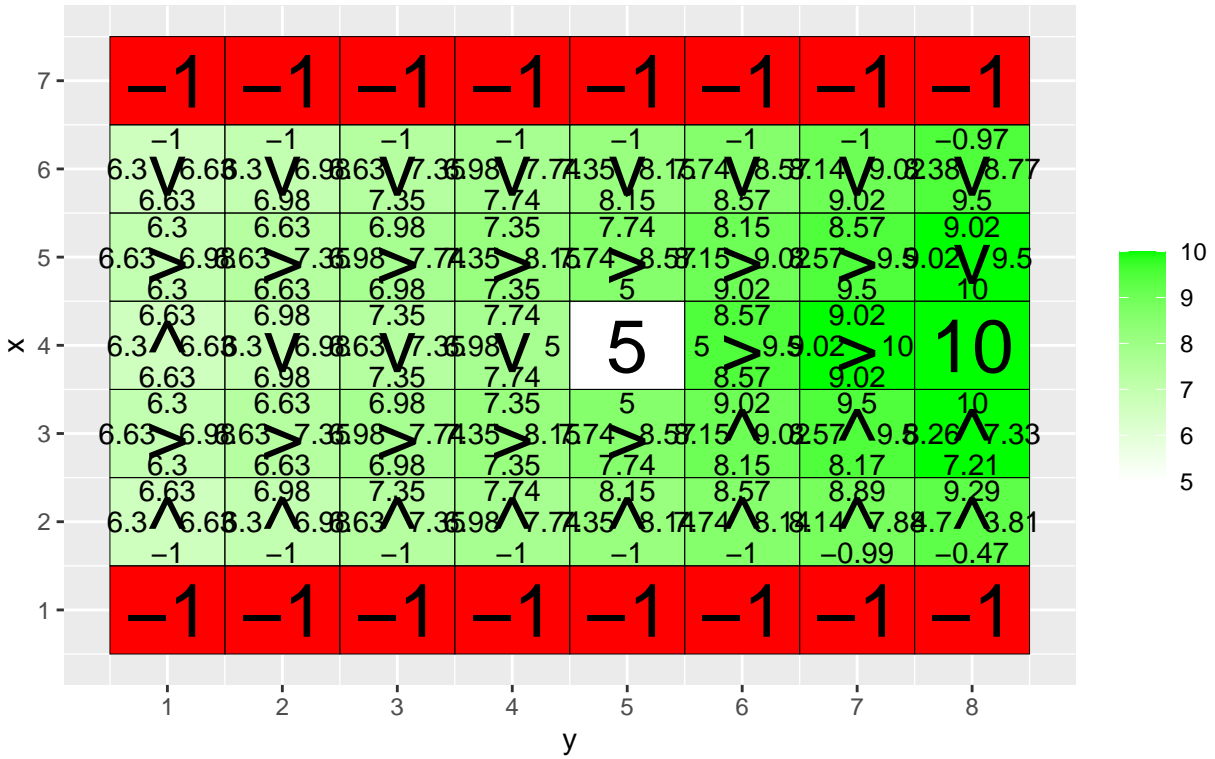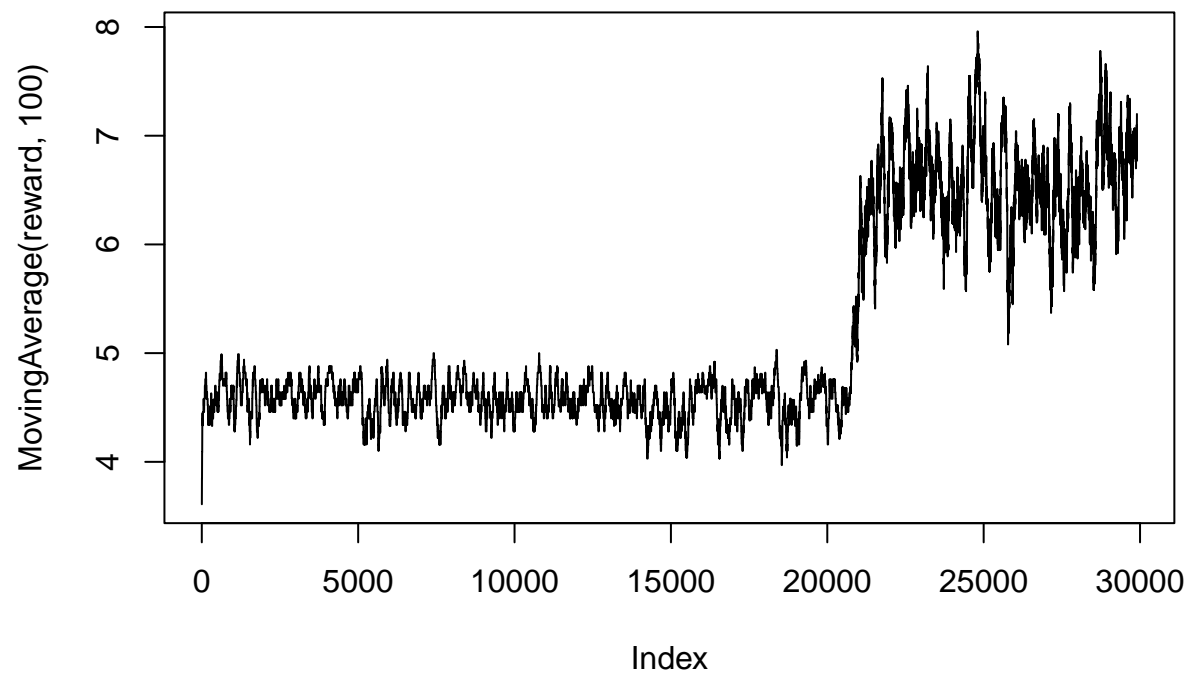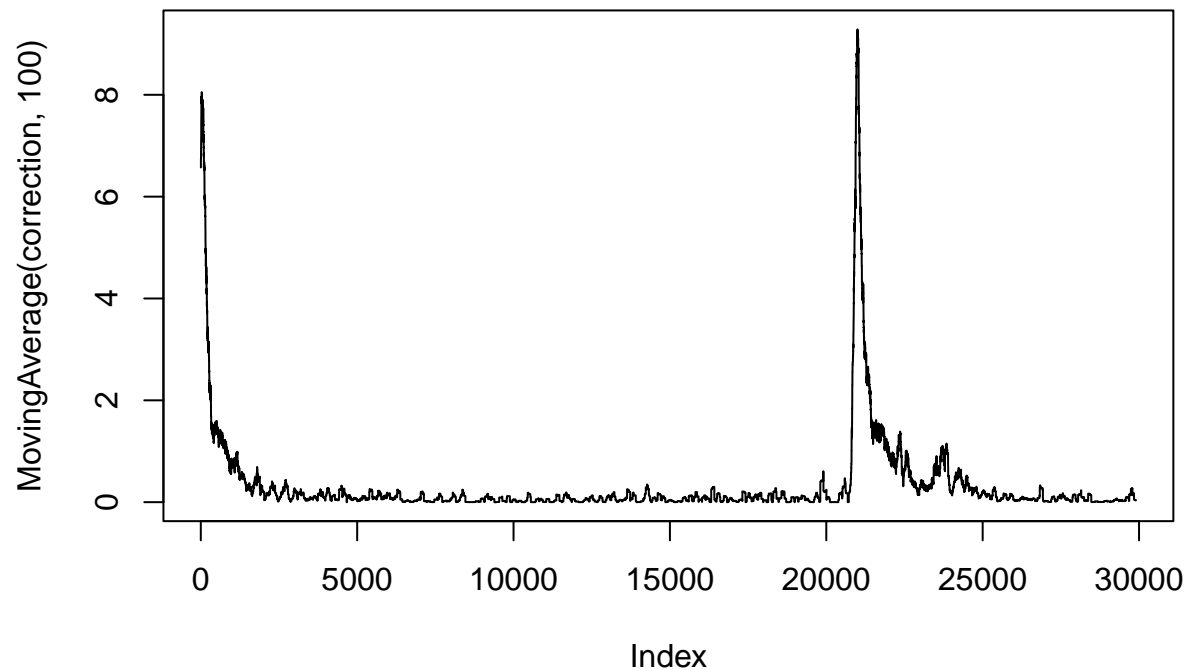(epsilon = 0.5 , alpha = 0.1 gamma = 0.75 , beta = 0 )

Q−table after 30000 iterations
(epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )

```
for(j in c(0.5,0.75,0.95)){
  q_table <- array(0,dim = c(H,W,4))
  reward <- NULL
  correction <- NULL

  for(i in 1:30000){
    # This is epsilon 0.1
    foo <- q_learning(epsilon = 0.1, gamma = j, start_state = c(4,1))
    reward <- c(reward,foo[1])
    correction <- c(correction,foo[2])
  }

  vis_environment(i, epsilon = 0.1, gamma = j)
  plot(MovingAverage(reward,100),type = "l")
  plot(MovingAverage(correction,100),type = "l")
}
```
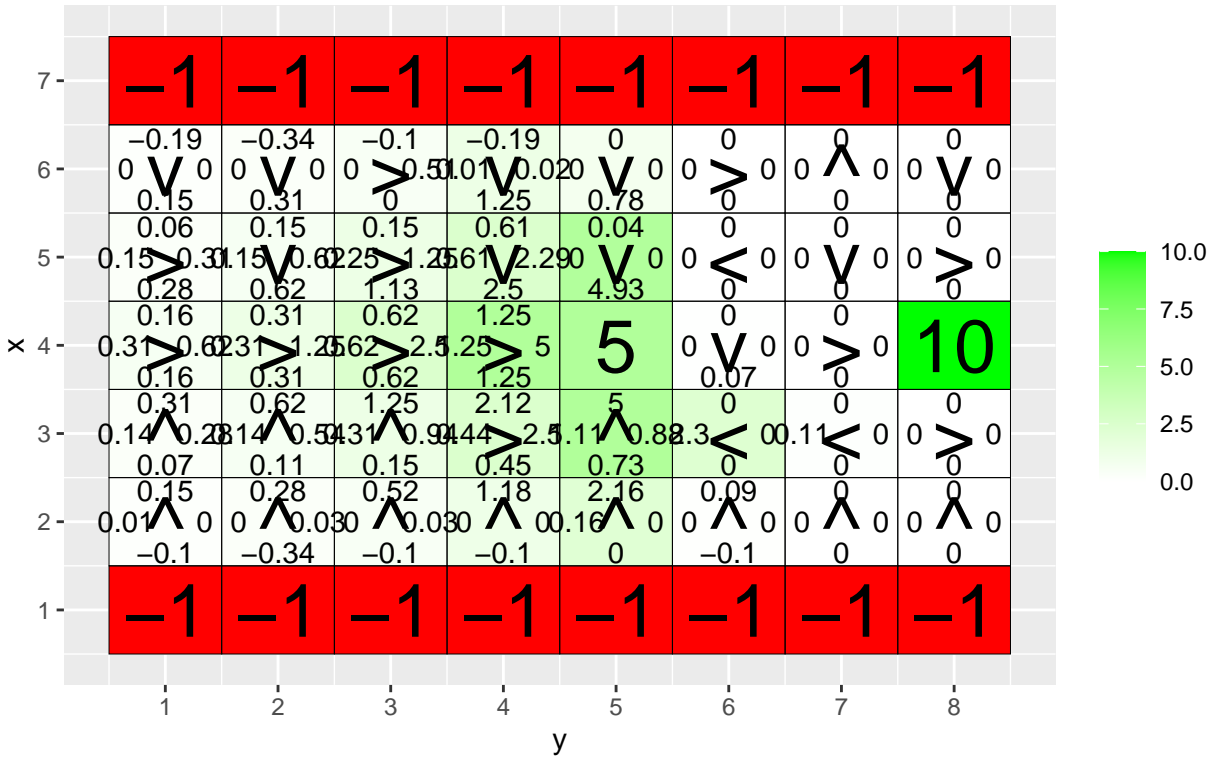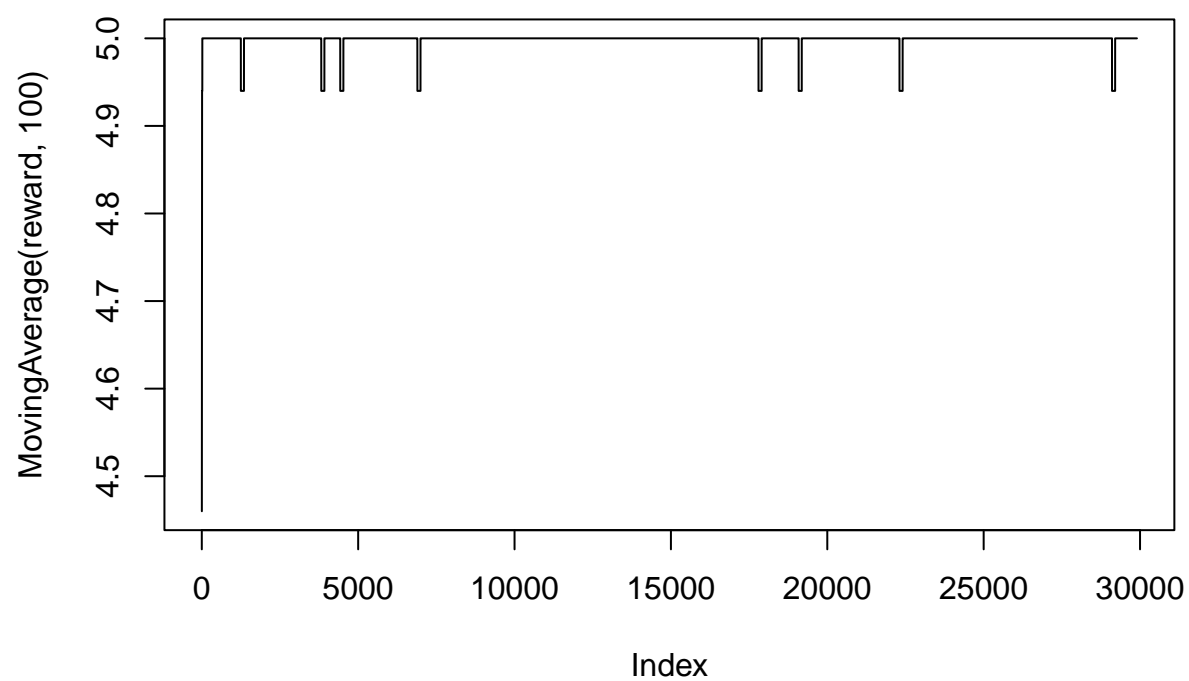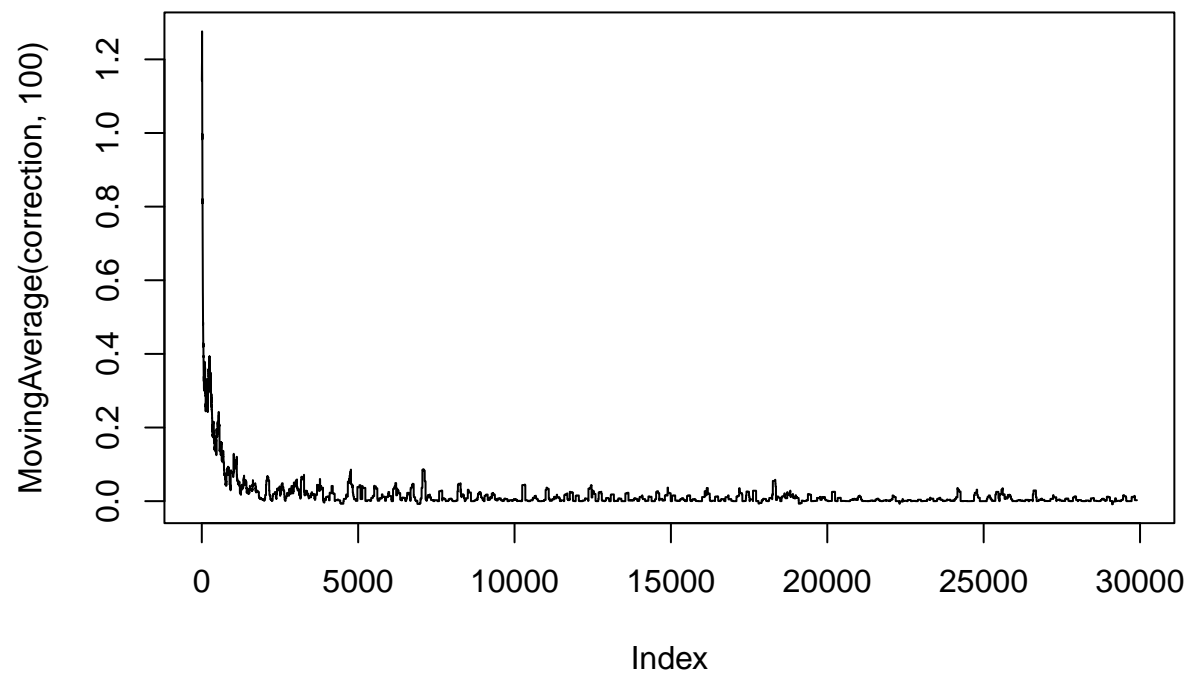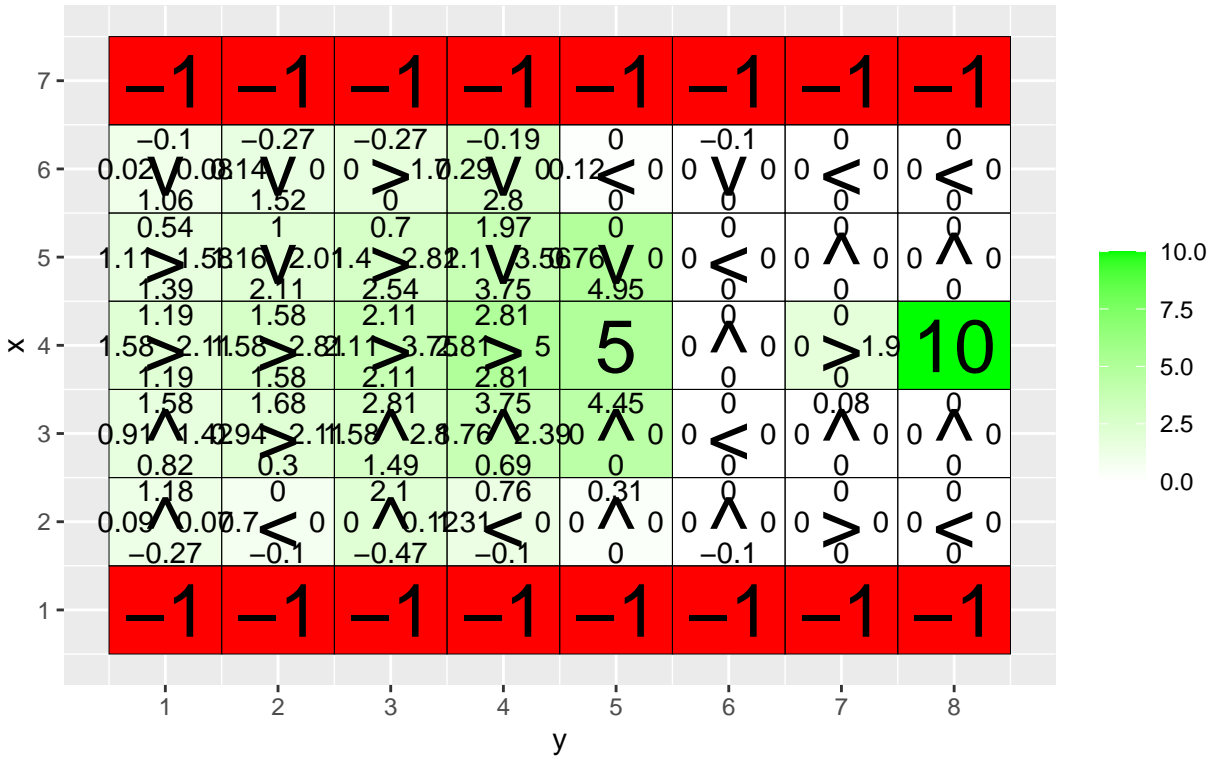
Q–table after 30000 iterations
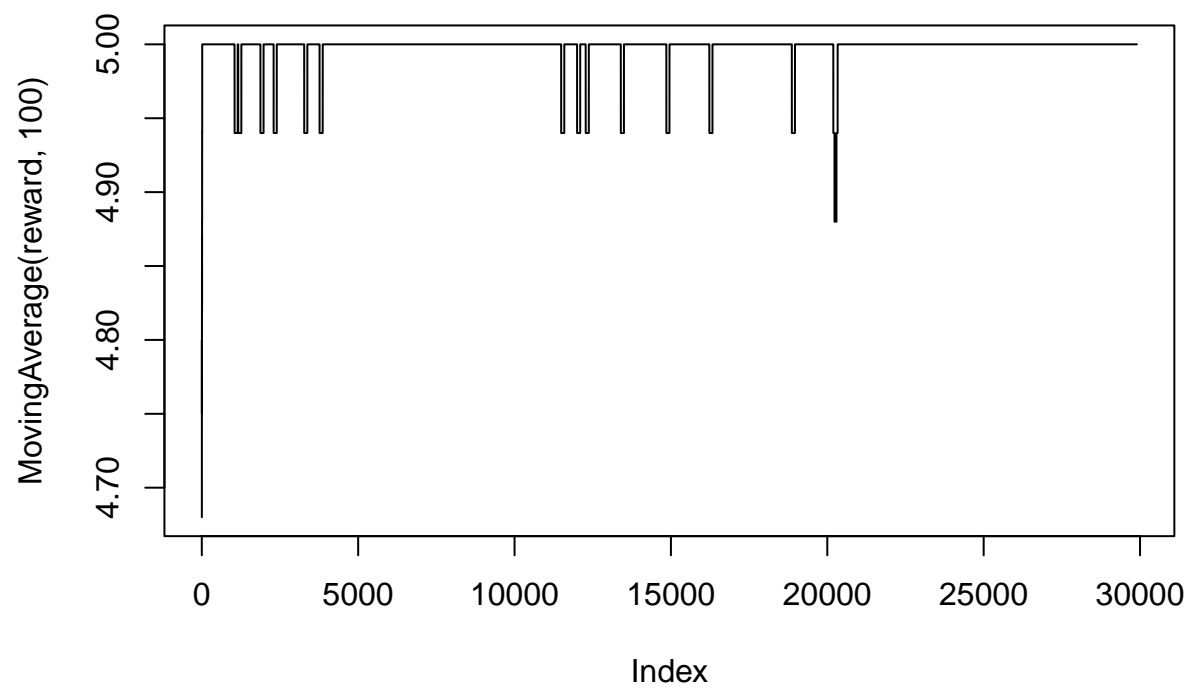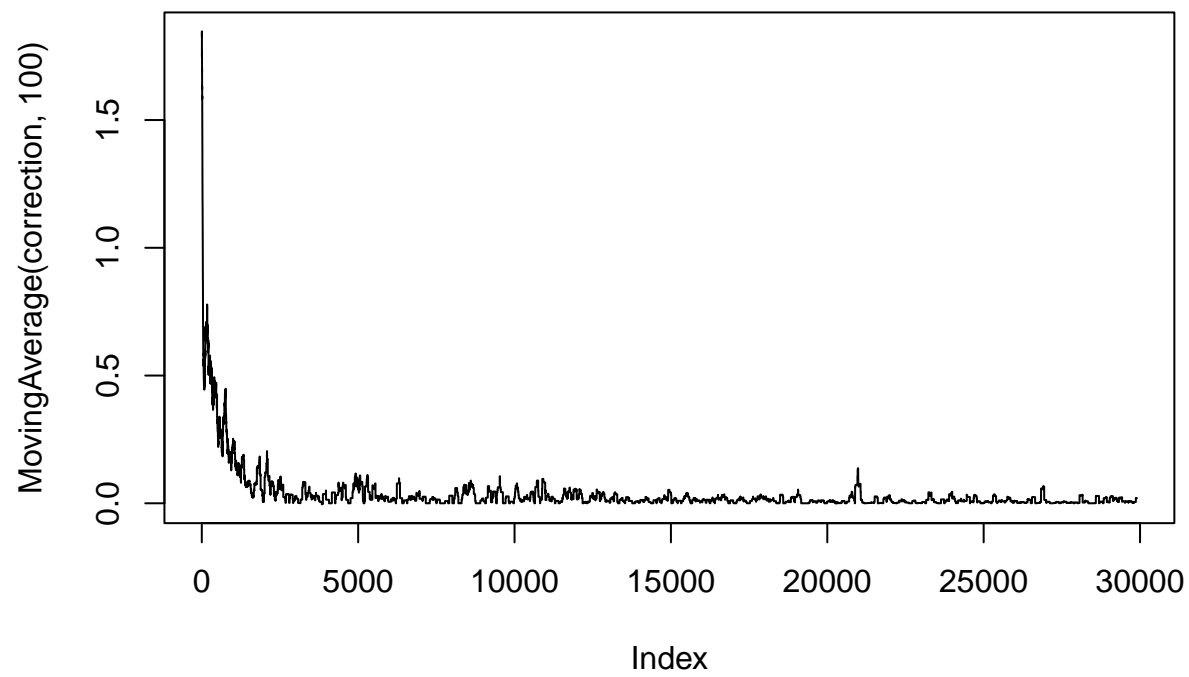(epsilon = 0.1 , alpha = 0.1 gamma = 0.5 , beta = 0 )

# Q–table after 30000 iterations
## (epsilon = 0.1 , alpha = 0.1 gamma = 0.75 , beta = 0 )

# Q−table after 30000 iterations
## (epsilon = 0.1 , alpha = 0.1 gamma = 0.95 , beta = 0 )

**Your task is to investigate how the epsilon and gamma parameters affect the learned policy by running**

My observations are:

**epsilon = 0.5**

**gamma = 0.5**: It chooses 5 quite often since the discount factor gamma is low meaning that it prioritizes future rewards low

**gamma = 0.75**: Most often chooses 10 but sometimes 5

**gamma = 0.95**: results in that it learns to avoid 5 reward and go for 10 reward

With epsilon = 0.1 the agent doesn't explore as much and subsequently doesn't find the 10 reward. It only finds 5 reward.

## 2.4 Environment C

Explore how the beta parameter affects the learned policy in a 3×6 environment by running 10,000 episodes of Q-learning with different beta values and analyzing the outcomes.

```
H <- 3
W <- 6

reward_map <- matrix(0, nrow = H, ncol = W)
reward_map[1,2:5] <- -1
reward_map[1,6] <- 10
```

```r
q_table <- array(0,dim = c(H,W,4))

vis_environment()
```

### Q–table after 0 iterations
### (epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )



```r
for(j in c(0,0.2,0.4,0.66)){
  q_table <- array(0,dim = c(H,W,4))

  for(i in 1:10000)
    foo <- q_learning(gamma = 0.6, beta = j, start_state = c(1,1))
    vis_environment(i, gamma = 0.6, beta = j)
}
```

# Q−table after 10000 iterations
## (epsilon = 0.5 , alpha = 0.1 gamma = 0.6 , beta = 0 )

# Q–table after 10000 iterations
## (epsilon = 0.5 , alpha = 0.1 gamma = 0.6 , beta = 0.2 )

Q–table after  10000  iterations
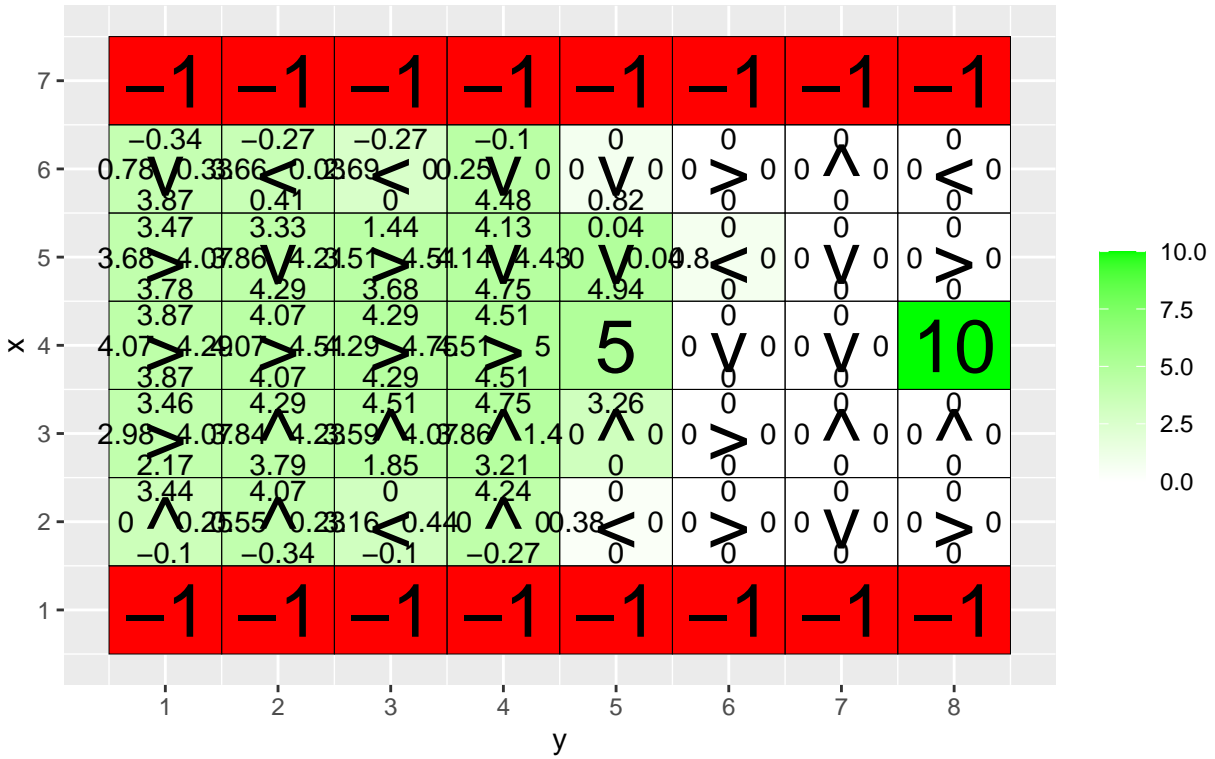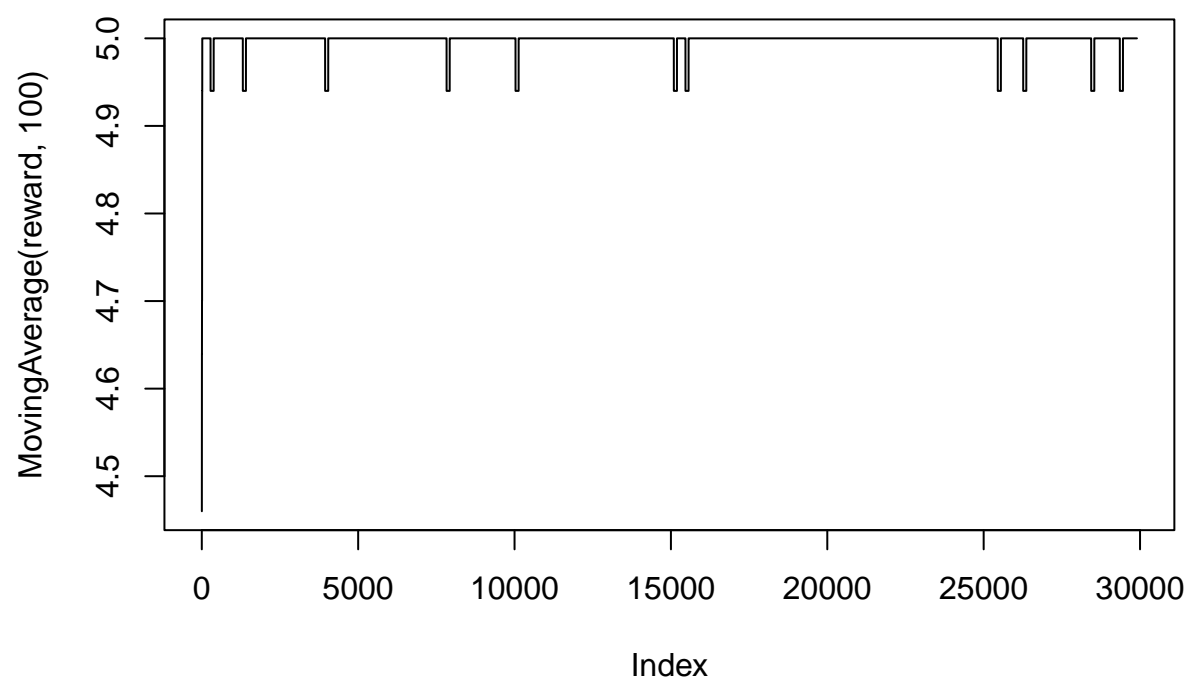(epsilon =  0.5 , alpha =  0.1 gamma =  0.6 , beta =  0.4 )
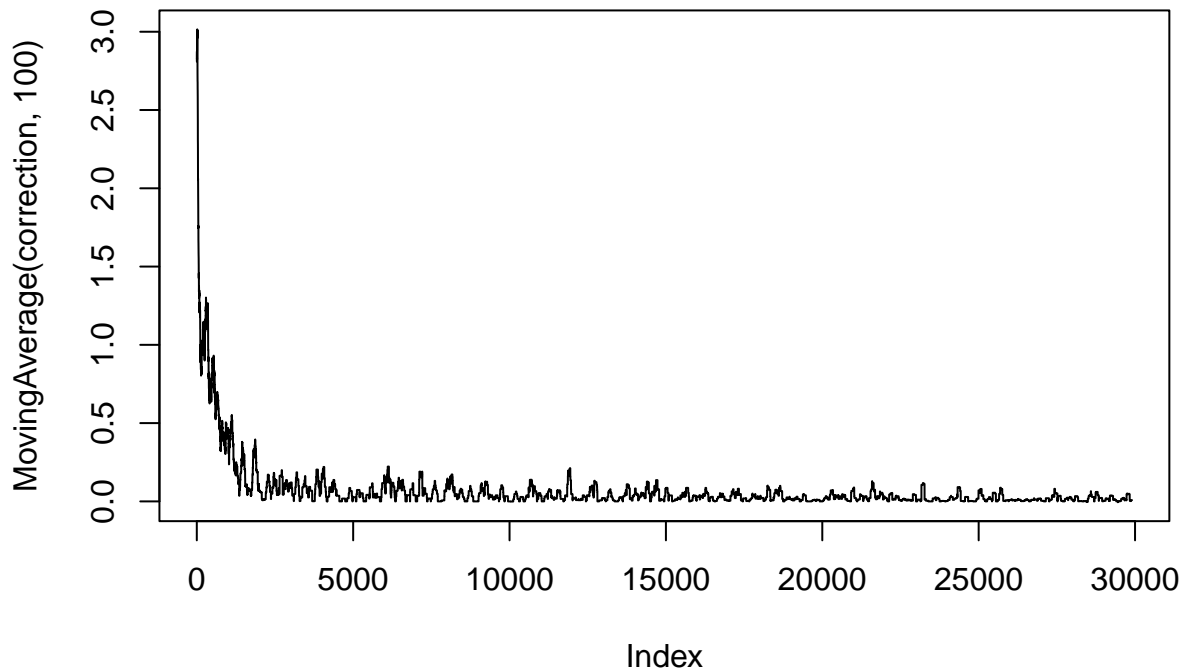
**Q–table after 10000 iterations**
**(epsilon = 0.5 , alpha = 0.1 gamma = 0.6 , beta = 0.66 )**

Here, Beta determines the probability that the agent will slip. A higher beta results in the agent avoiding the -1 wall with greater margin since it might slip.

## 2.5 REINFORCE

Study the provided REINFORCE algorithm for a 4×4 grid-world, where the agent must learn to navigate to a random goal position, and analyze the code and results without needing to run it.

Plot 1 and 2

## Action probabilities after 5000 episodes



## Action probabilities after 5000 episodes

Plot 3 and 4

Action probabilities after 5000 episodes

Row 4:
- Cell (1,4): top 0, left 0, center V, right 0.01, bottom 0.99
- Cell (2,4): top 0, left 0, center V, right 0.01, bottom 0.99
- Cell (3,4): top 0, left 0.05, center V, right 0.01, bottom 0.95
- Cell (4,4): top 0, left 0.35, center V, right 0.01, bottom 0.64

Row 3:
- Cell (1,3): top 0, left 0, center V, right 0.02, bottom 0.97
- Cell (2,3): top 0, left 0.04, center V, right 0.02, bottom 0.94
- Cell (3,3): top 0, left 0.32, center V, right 0.02, bottom 0.65
- Cell (4,3): top 0, left 0.86, center <, right 0.01, bottom 0.13

Row 2:
- Cell (1,2): top 0, left 0.03, center V, right 0.08, bottom 0.88
- Cell (2,2): top 0, left 0.29, center V, right 0.07, bottom 0.64
- Cell (3,2): top 0, left 0.84, center <, right 0.02, bottom 0.14
- Cell (4,2): top 0, left 0.98, center <, right 0, bottom 0.01

Row 1:
- Cell (1,1): Goal
- Cell (2,1): top 0.07, left 0.77, center <, right 0.07, bottom 0.08
- Cell (3,1): top 0.01, left 0.95, center <, right 0.01, bottom 0.02
- Cell (4,1): top 0, left 0.98, center <, right 0, bottom 0.01

x-axis: 1  2  3  4
y

Action probabilities after 5000 episodes

Row 4:
- Cell (1,4): top 0, left 0, center V, right 0.01, bottom 0.99
- Cell (2,4): top 0, left 0, center V, right 0, bottom 1
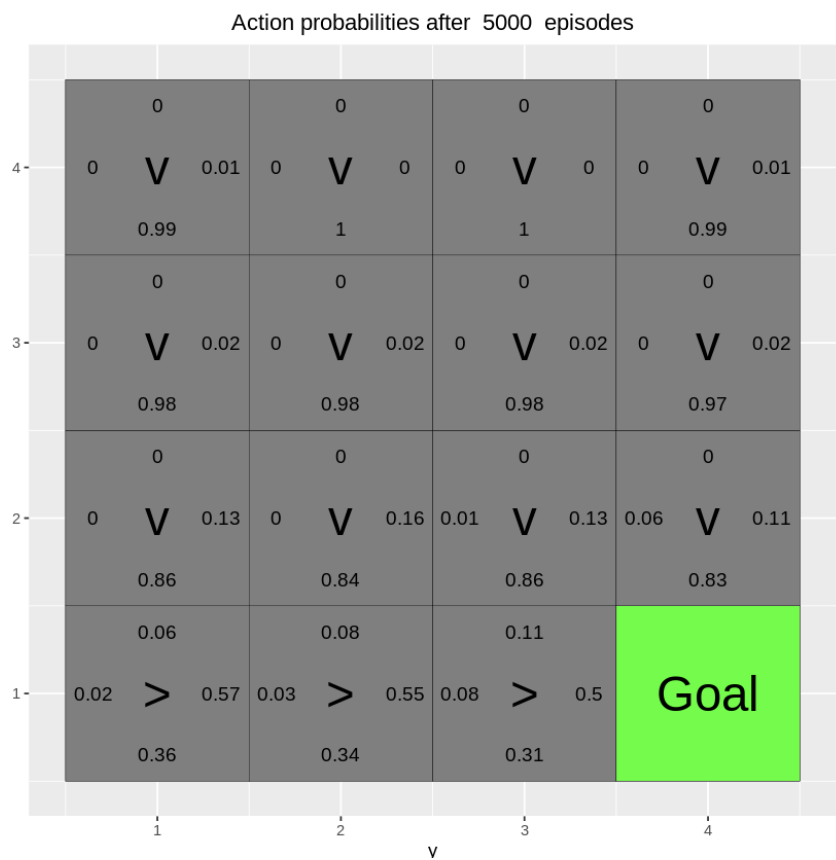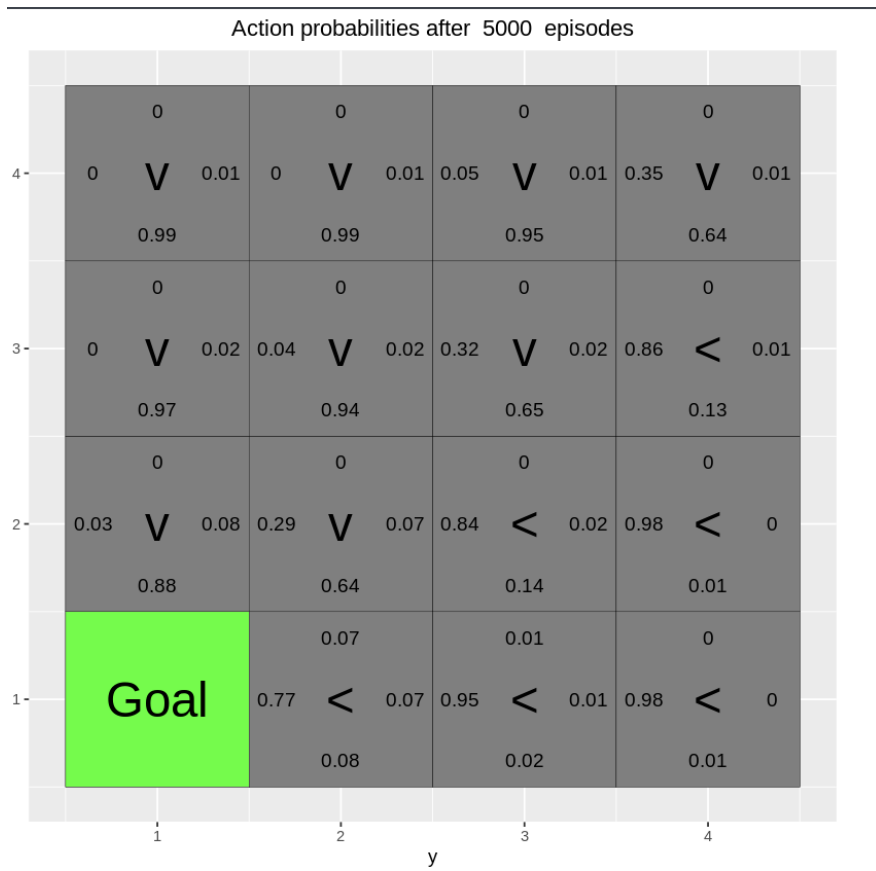- Cell (3,4): top 0, left 0, center V, right 0, bottom 1
- Cell (4,4): top 0, left 0, center V, right 0.01, bottom 0.99

Row 3:
- Cell (1,3): top 0, left 0, center V, right 0.02, bottom 0.98
- Cell (2,3): top 0, left 0, center V, right 0.02, bottom 0.98
- Cell (3,3): top 0, left 0, center V, right 0.02, bottom 0.98
- Cell (4,3): top 0, left 0, center V, right 0.02, bottom 0.97

Row 2:
- Cell (1,2): top 0, left 0, center V, right 0.13, bottom 0.86
- Cell (2,2): top 0, left 0, center V, right 0.16, bottom 0.84
- Cell (3,2): top 0, left 0.01, center V, right 0.13, bottom 0.86
- Cell (4,2): top 0, left 0.06, center V, right 0.11, bottom 0.83

Row 1:
- Cell (1,1): top 0.06, left 0.02, center >, right 0.57, bottom 0.36
- Cell (2,1): top 0.08, left 0.03, center >, right 0.55, bottom 0.34
- Cell (3,1): top 0.11, left 0.08, center >, right 0.5, bottom 0.31
- Cell (4,1): Goal

x-axis: 1  2  3  4
v

## 2.6 Environment D

Train a REINFORCE agent on eight goal positions and validate it on the remaining eight positions, then answer questions on whether the agent learned a good policy and whether Q-learning could solve the task.

**Has the agent learned a good policy? Why / Why not ?**

Yes, the agent has learned relatively good policy (plot 3 and 4). The policy generalizes well during validation since the training goals were spread across the grid, allowing the agent to learn a flexible policy.

**Could you have used the Q-learning algorithm to solve this task?**

Yes, although REINFORCE is much more suitable. Q-learning could solve this task by incorporating the goal coordinates into the state representation. However, Q-learning might struggle with generalizing as effectively as REINFORCE.

## 2.7 Environment E

Analyze the agent's performance when trained with goals from the top row and validated with goals from lower rows, and compare the results with environment D, explaining any differences.

**Has the agent learned a good policy? Why / Why not ?**

No, the agent has not learned a fully generalized policy (plot 1 and 2). It has overfitted to the training where the goals were located in the top row making the policy arrows pointing upward.

**If the results obtained for environments D and E differ, explain why.**

The difference arises from the diversity of the training goals. In Environment D, the training goals were spread across the grid, resulting in better generalization. In Environment E, the training goals were restricted to the top row, leading to overfitting and poor generalization to new goal positions during validation.