# TDDE15_Lab3_oskhi827

Oskar Hidén - oskhi827

10/5/2020

```r
require(extrafont)
    # need only do this once!
    font_import(pattern="[A/a]rial", prompt=FALSE)
# By Jose M. Peña and Joel Oskarsson.
# For teaching purposes.
# jose.m.pena@liu.se.


#############################################################################################
# Q-learning
#############################################################################################

# install.packages("ggplot2")
# install.packages("vctrs")
library(ggplot2)

# If you do not see four arrows in line 16, then do the following:
# File/Reopen with Encoding/UTF-8

arrows <- c("↑", "→", "↓", "←")
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){

  # Visualize an environment with rewards.
  # Q-values for all actions are displayed on the edges of each tile.
  # The (greedy) policy for each state is also displayed.
  #
  # Args:
  #   iterations, epsilon, alpha, gamma, beta (optional): for the figure title.
  #   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.
  #   H, W (global variables): environment dimensions.

  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)
```

```r
    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) +
            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 a greedy action for state (x,y) from q_table.
  #
  # Args:
  #   x, y: state coordinates.
  #   q_table (global variable): a HxWx4 array containing Q-values for each state-action pair.
  #
  # Returns:
  #   An action, i.e. integer in {1,2,3,4}.

  # Your code here.
  return(which.max(rank(q_table[x,y,], ties.method = "random")))
}

EpsilonGreedyPolicy <- function(x, y, epsilon){

  # Get an epsilon-greedy action for state (x,y) from q_table.
  #
  # Args:
  #   x, y: state coordinates.
  #   epsilon: probability of acting randomly.
  #
  # Returns:
  #   An action, i.e. integer in {1,2,3,4}.
```

```r
  # Your code here.
  if(runif(1)>epsilon){
    direction = GreedyPolicy(x,y)
  }else{
    direction = sample(1:4,1)
  }

  return(direction)
}

transition_model <- function(x, y, action, beta){

  # Computes the new state after given action is taken. The agent will follow the action
  # with probability (1-beta) and slip to the right or left with probability beta/2 each.
  #
  # Args:
  #   x, y: state coordinates.
  #   action: which action the agent takes (in {1,2,3,4}).
  #   beta: probability of the agent slipping to the side when trying to move.
  #   H, W (global variables): environment dimensions.
  #
  # Returns:
  #   The new state after the action has been taken.

  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.
  #
  # Args:
  #   start_state: array with two entries, describing the starting position of the agent.
  #   epsilon (optional): probability of acting greedily.
  #   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 <<-.
```

```r
  # Your code here.
  current_state = start_state
  episode_correction = 0
  repeat{
    # Follow policy, execute action, get reward.
    action = EpsilonGreedyPolicy(current_state[1], current_state[2], epsilon)
    new_state = transition_model(current_state[1], current_state[2], action, beta)
    reward = reward_map[new_state[1], new_state[2]]

    # Q-table update.
    q_action_value = q_table[current_state[1], current_state[2], action]
    next_exp_r = max(q_table[new_state[1], new_state[2],])

    correction = gamma*next_exp_r - q_action_value
    q_table[current_state[1], current_state[2],action] <<- q_action_value + alpha*(reward + correction)
    episode_correction = episode_correction+correction
    current_state=new_state
    if(reward!=0)
      # End episode.
      return (c(reward,episode_correction))

  }

}

#######################################################################################
# Q-Learning Environments
#######################################################################################

# Environment A (learning)

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 )

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

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| **5** | 0 / 0 ... 0 / 0 | 0 / 0 ... 0 / 0 | 0 / 0 ... 0 / −0.27 | 0 / 0 ... 0 / 0 | 0 / 0 ... 0 / 0 | 0 / 0 ... 0 / 0 | 0 / 0 ... 0 / 0 |
| **4** | 0 / 0 ... 0 / 0.01 | 0 / 0 −0.61 ... 0 / 0.03 | **−1** | 0 / −0.19 0.27 ... 0 / 0 | 0 / 0 ... 0 / 2.52 | 0 / 0 ... 0 / 1 | 0 / 0 ... 0 / 0 |
| **3** | 0 / 0.02 0.11 ... 0 / 0.03 | 0 / 0 −0.85 ... 0 / 0.22 | **−1** | 0.01 / −0.19 6.53 ... 0 / 0 | 0.34 / 3.61 9.28 ... 0 / 0 | **10** | 0 / 0 ... 0 / 0 |
| **2** | 0 / 0 0.12 ... 0 / 0.01 | 0 / 2.02 −0.92 ... 0 / 0.43 | **−1** | 3.82 / −0.34 0.88 ... 0 / 0.41 | 2.85 / 8.36 0 ... 0 / 0.01 | 1 / 0 0 ... 0 / 0 | 0 / 0 ... 0 / 0 |
| **1** | 0.01 / 0 0.12 ... 0 / 0 | 0.04 / 4.01 0.87 ... 0 / 0.13 | −0.69 / 7.11 1.49 ... 0 / 0.24 | 2.61 / 0.3 0.03 ... 0 / 0.35 | 0 / 3.5 0 ... 0 / 0 | 0 / 0 0 ... 0 / 0 | 0 / 0 ... 0 / 0 |

x (vertical axis)

y (horizontal axis)

Legend: 10.0 / 7.5 / 5.0 / 2.5 / 0.0

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

x

Row 5:
| 2.64 | 1.88 | 0.01 | 0.11 | 0 | 0.07 | 0 |
1.53  2.4 | .13  0.00 | 7.02  1.1 | 0.08  0.1 | 7.85  0.0 | 4.11  0 | 0.11  0
| 5.77 | 5.8 | −0.34 | 4.94 | 2.38 | 1.02 | 0 |

Row 4:
| 4.93 | 4.09 | | 1.37 | 0.24 | 0.17 | 0 |
5.86  65.33  −0.98 | **−1** | −0.72  5.5 | 9.84  2.7 | 0.96  0.0 | 6.41  0
| 6.3 | 6.63 | | 8.91 | 9.44 | 6.13 | 0 |

Row 3:
| 5.99 | 6.3 | | 8.11 | 8.77 | | 0 |
6.3  6.63 | 6.3  −1 | **−1** | −1  9.5 | 8.98  10 | **10**  1.9 | 0
| 6.63 | 6.98 | | 8.53 | 8.96 | | 0 |

Row 2:
| 6.21 | 6.63 | | 9.02 | 9.5 | 9.2 | 0.1 |
6.36  6.9 | 6.63  −1 | **−1** | −1  9.0 | 7.23  7.0 | 3.13  0.7 | 2.73  0.21
| 6.87 | 7.35 | | 8.12 | 6.05 | 0.59 | 0.06 |

Row 1:
| 6.3 | 6.98 | −1 | 8.57 | 9.02 | 1.42 | 0.58 |
5.33  7.35 | 6.98  7.74 | 7.35  8.1 | 3.74  8.5 | 6.54  3.5 | 3.57  0.0 | 0.05  0
| 5.87 | 7.35 | 7.74 | 8.14 | 7.11 | 1.25 | 0 |

y: 1  2  3  4  5  6  7

Legend:
10.0
7.5
5.0
2.5

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

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| **5** | 6.9 / 6.85 7.36 / 6.53 | 7.35 / 6.98 7.74 / 6.98 | 7.74 / 4.35 8.13 / −1 | 8.13 / 5.73 8.57 / 8.57 | 7.73 / 7.79 8.24 / 9.02 | 4.72 / 6.32 1.62 / 9.41 | 0.18 / 7.32 0.36 / 5 |
| **4** | 6.98 / 6.63 6.96 / 6.3 | 7.35 / 6.63 −1 / 6.63 | −1 | 8.14 / −1 9.02 / 9.02 | 8.57 / 8.57 9.57 / 9.5 | 7.63 / 7.53 5.78 / 10 | 0.85 / 8.94 3.15 / 2.92 |
| **3** | 6.63 / 6.3 6.63 / 6.63 | 6.98 / 6.3 −1 / 6.98 | −1 | 8.57 / −1 8.57 / 8.57 | 9.02 / 9.5 9.02 / 9.02 | 10 / 10 9.62 / 1.2 | 3.95 / 2.8 |
| **2** | 6.3 / 6.63 6.96 / 6.98 | 6.63 / 6.63 −1 / 7.35 | −1 | 9.02 / −1 9.02 / 8.15 | 9.5 / 8.57 9.5 / 8.57 | 10 / 9.02 9.02 / 9.02 | 8.18 / 9.5 8.99 / 8.55 |
| **1** | 6.63 / 6.98 7.36 / 6.98 | 6.98 / 6.98 7.74 / 7.35 | −1 / 4.35 8.13 / 7.74 | 8.57 / 5.74 8.57 / 8.15 | 9.02 / 7.15 9.02 / 8.57 | 9.5 / 8.57 8.58 / 9.02 | 9.02 / 8.01 8.15 / 7.57 |

x (vertical axis) · y (horizontal axis: 1 2 3 4 5 6 7)

Legend: 10 9 8 7 6 5

```r
# Environment B (the effect of epsilon and gamma)

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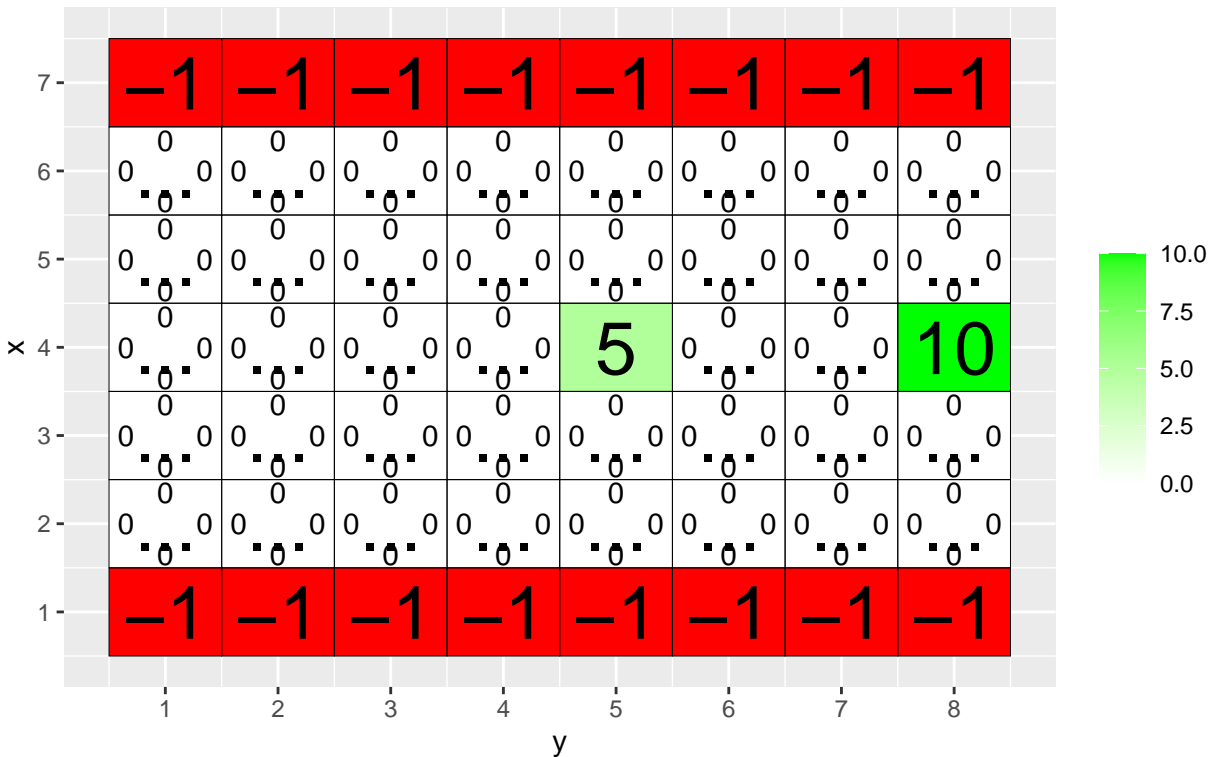
Q–table after  0  iterations
(epsilon =  0.5 , alpha =  0.1 gamma =  0.95 , beta =  0 )

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

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

  return (rsum)
}

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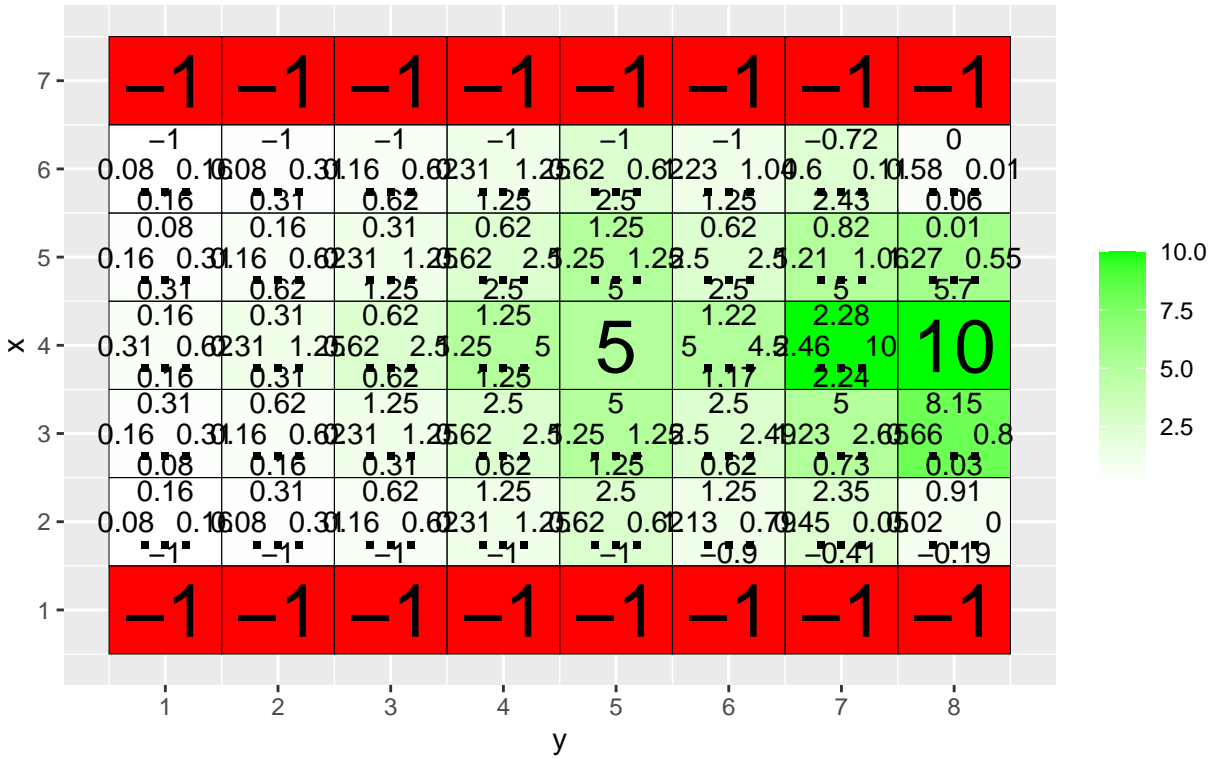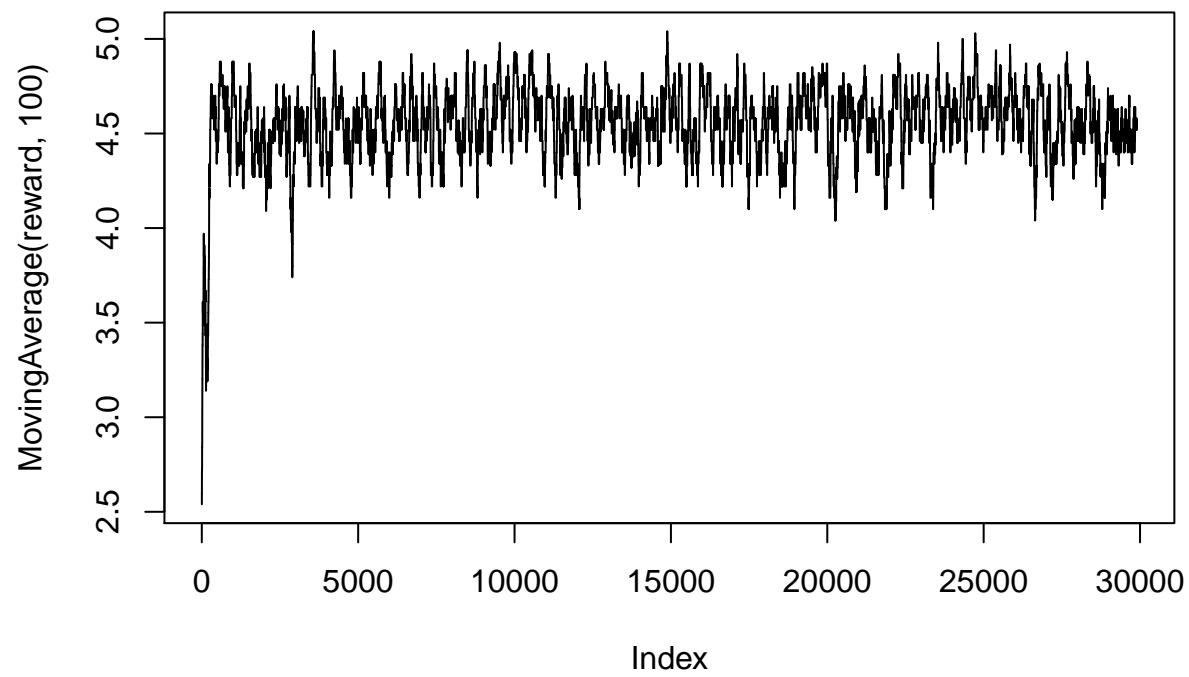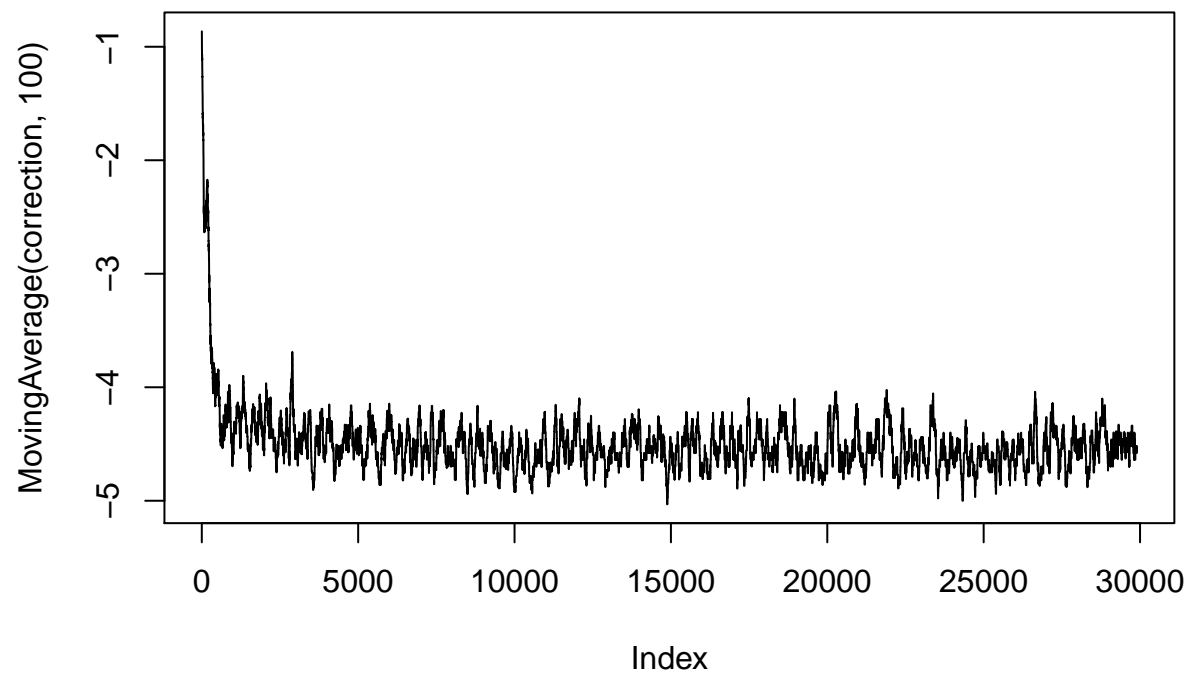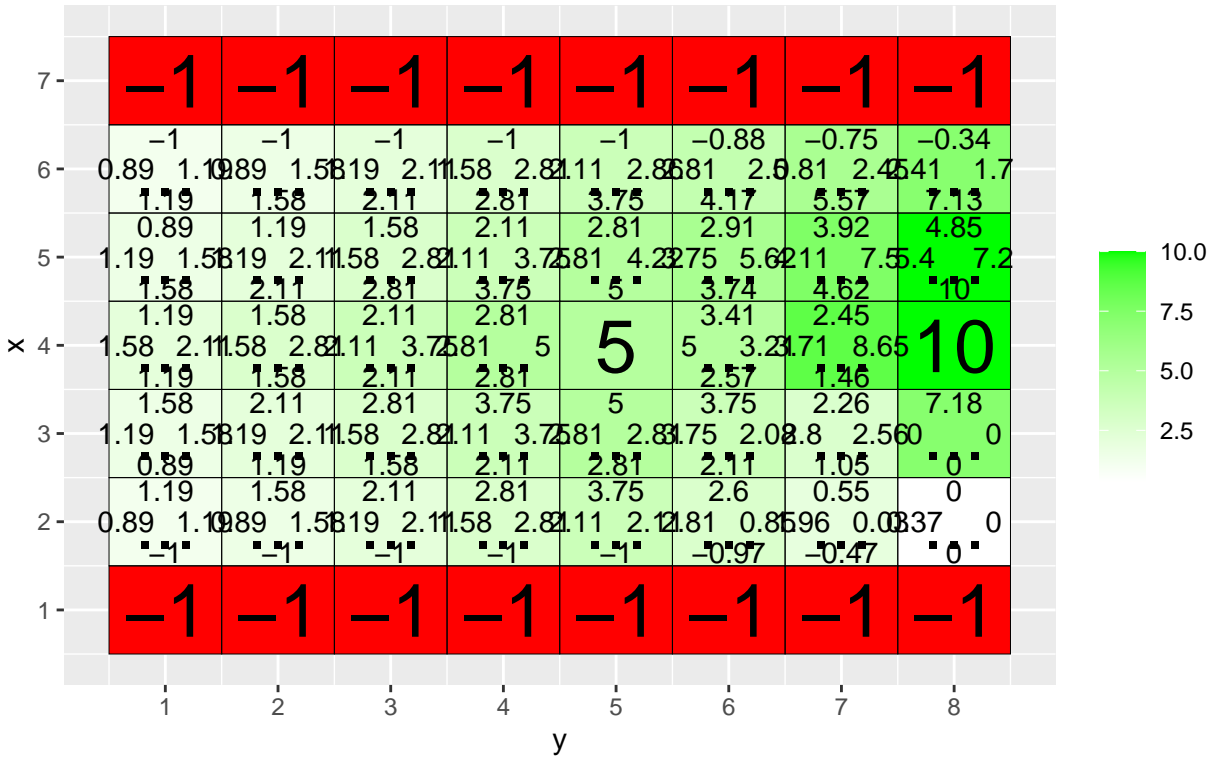
Q−table after 30000 iterations
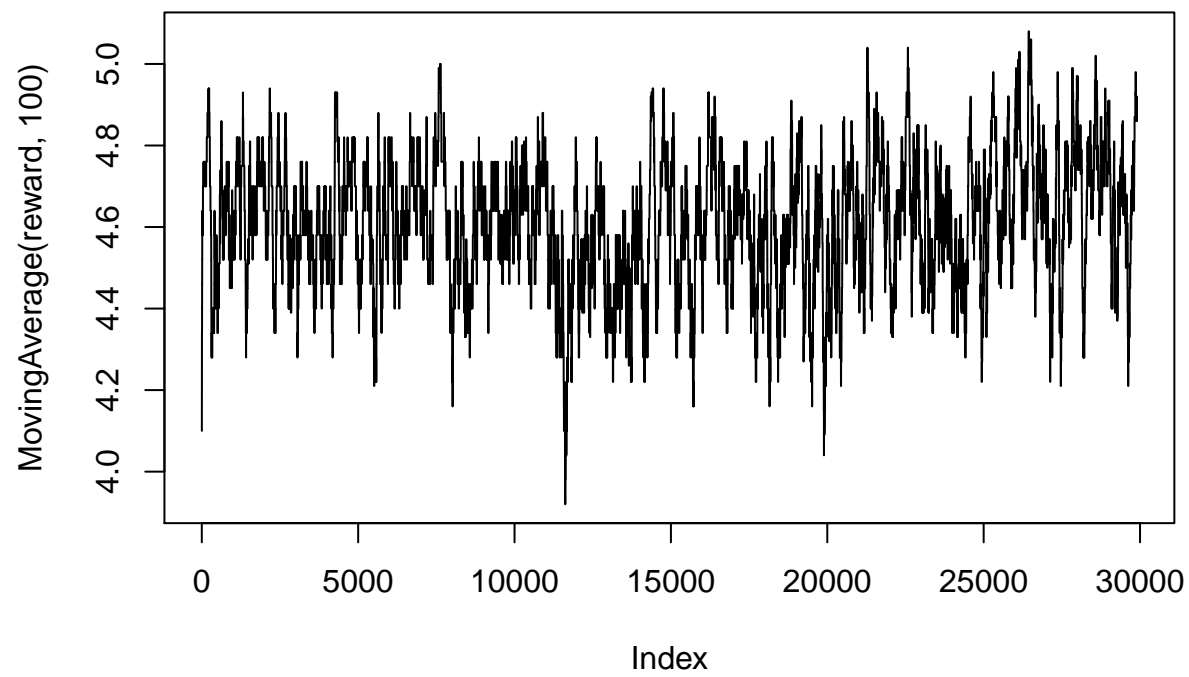(epsilon = 0.5 , alpha = 0.1 gamma = 0.5 , beta = 0 )

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

x

−1 −1 −1 −1 −1 −1 −1 −1

−1      −1      −1      −1      −1      −0.88   −0.75   −0.34
7 -
6 - 0.89  1.10.89  1.5819  2.11.58  2.81.11  2.8581  2.6.81  2.4541   1.7
1.19    1.58    2.11    2.81    3.75    4.17    5.57    7.13

0.89    1.19    1.58    2.11    2.81    2.91    3.92    4.85
5 - 1.19  1.5819  2.11.58  2.81.11  3.7581  4.2375  5.6211   7.5.4   7.2
1.58    2.11    2.81    3.75     5      3.74    4.62    10

1.19    1.58    2.11    2.81            3.41    2.45
4 - 1.58  2.11.58  2.81.11  3.7581   5     5    3.23.71  8.65    10
1.19    1.58    2.11    2.81     5      2.57    1.46

1.58    2.11    2.81    3.75     5      3.75    2.26    7.18
3 - 1.19  1.5819  2.11.58  2.81.11  3.7581  2.83.75  2.02.8  2.560    0
0.89    1.19    1.58    2.11    2.81    2.11    1.05     0

1.19    1.58    2.11    2.81    3.75    2.6     0.55     0
2 - 0.89  1.10.89  1.5819  2.11.58  2.81.11  2.181  0.8596  0.0337   0
−1      −1      −1      −1      −1      −0.97   −0.47    0

−1 −1 −1 −1 −1 −1 −1 −1
1 -

        1       2       3       4       5       6       7       8
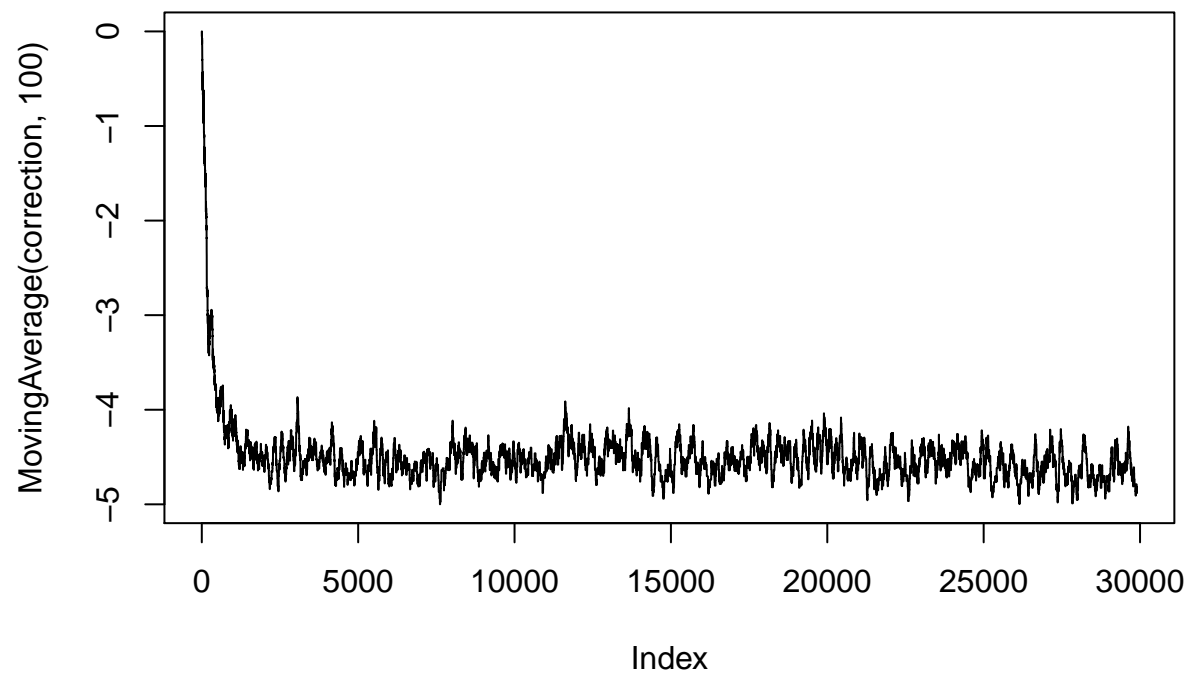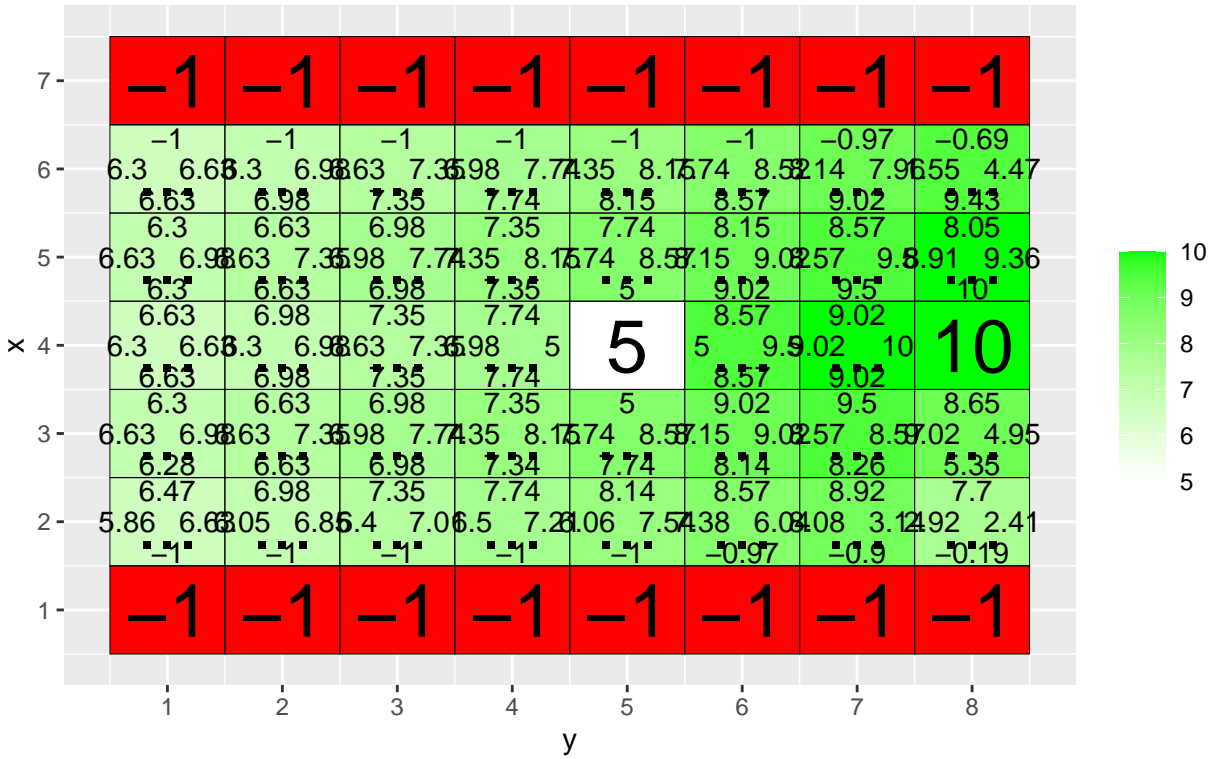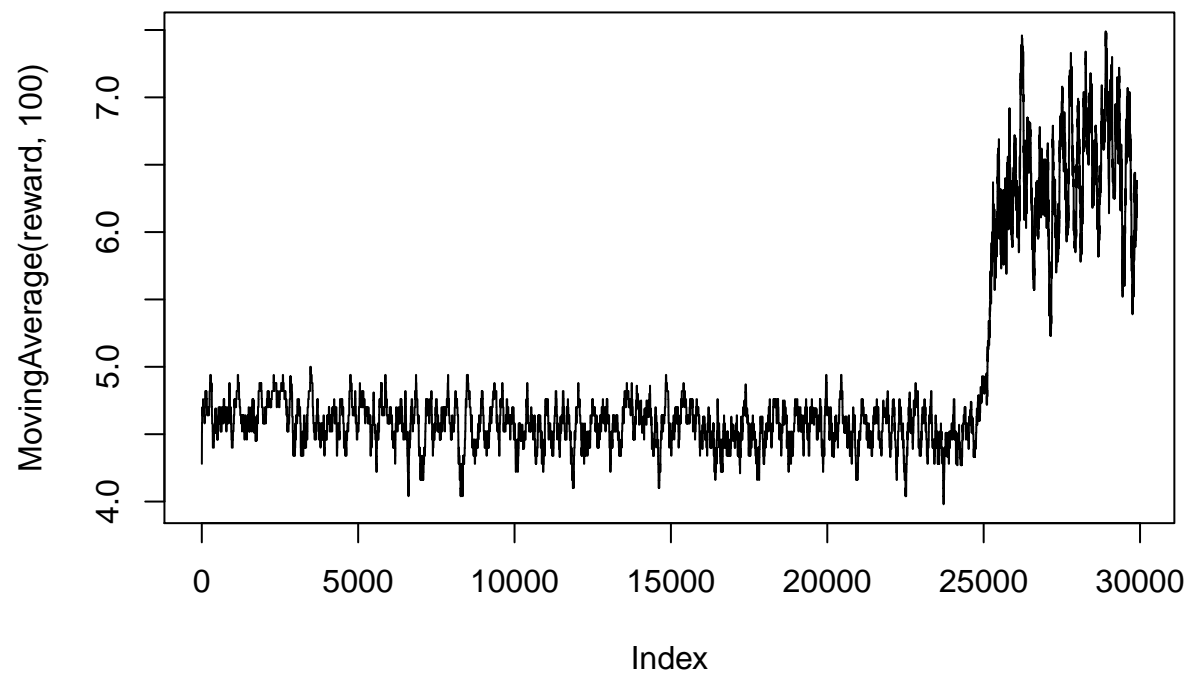
y

10.0

7.5

5.0

2.5

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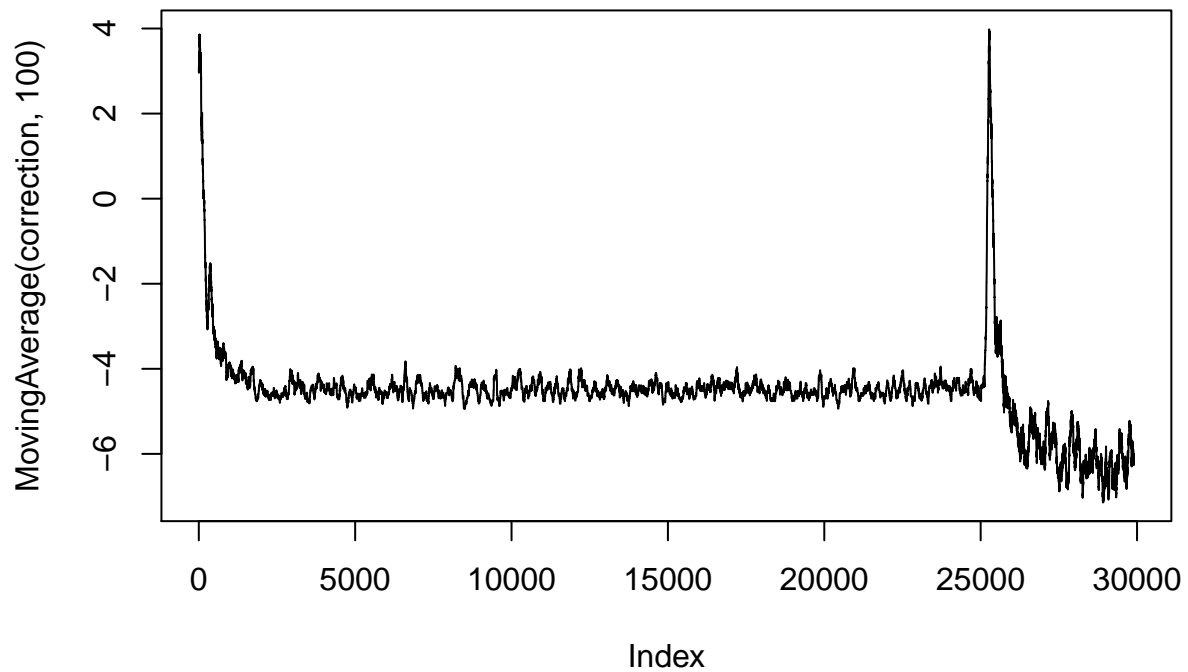(epsilon =  0.5 , alpha =  0.1 gamma =  0.95 , beta =  0 )

x

7 - | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1

−1  −1  −1  −1  −1  −1  −0.97  −0.69
6 - 6.3  6.6 6.3  6.9 6.63  7.35 6.98  7.74 7.35  8.15 7.74  8.57 8.14  7.96 6.55  4.47
6.63  6.98  7.35  7.74  8.15  8.57  9.02  9.43

6.3  6.63  6.98  7.35  7.74  8.15  8.57  8.05
5 - 6.63  6.9 6.63  7.35 6.98  7.74 7.35  8.15 7.74  8.57 8.15  9.02 8.57  9.8 8.91  9.36
6.3  6.63  6.98  7.35  5  9.02  9.5  10

6.63  6.98  7.35  7.74       8.57  9.02
4 - 6.3  6.6 6.3  6.9 6.63  7.35 6.98  5 | 5 | 5  9.9 9.02  10 | 10
6.63  6.98  7.35  7.74       8.57  9.02

6.3  6.63  6.98  7.35  5  9.02  9.5  8.65
3 - 6.63  6.9 6.63  7.35 6.98  7.74 7.35  8.15 7.74  8.57 8.15  9.02 8.57  8.5 7.02  4.95
6.28  6.63  6.98  7.34  7.74  8.14  8.26  5.35

6.47  6.98  7.35  7.74  8.14  8.57  8.92  7.7
2 - 5.86  6.6 6.05  6.8 6.4  7.06 6.5  7.2 7.06  7.54 7.38  6.08 4.08  3.12 4.92  2.41
−1  −1  −1  −1  −1  −0.97  −0.9  −0.19

1 - | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1

       1       2       3       4       5       6       7       8

y

10
9
8
7
6
5

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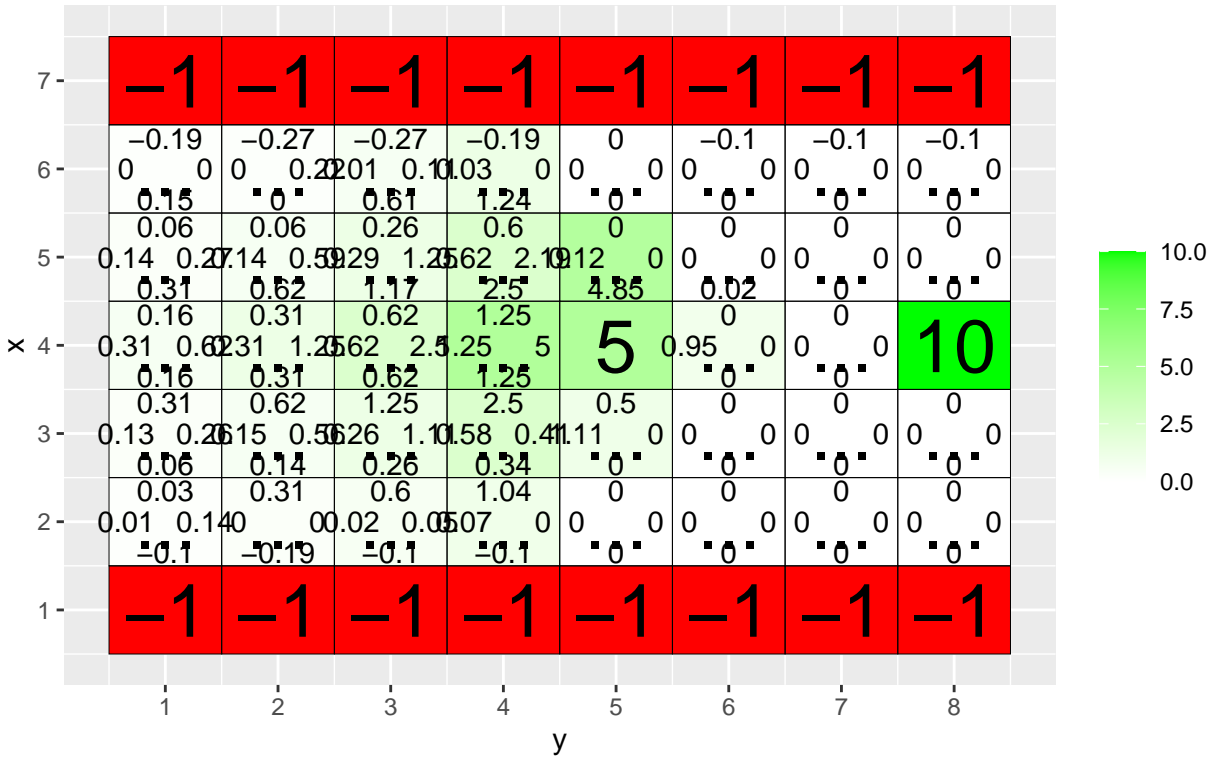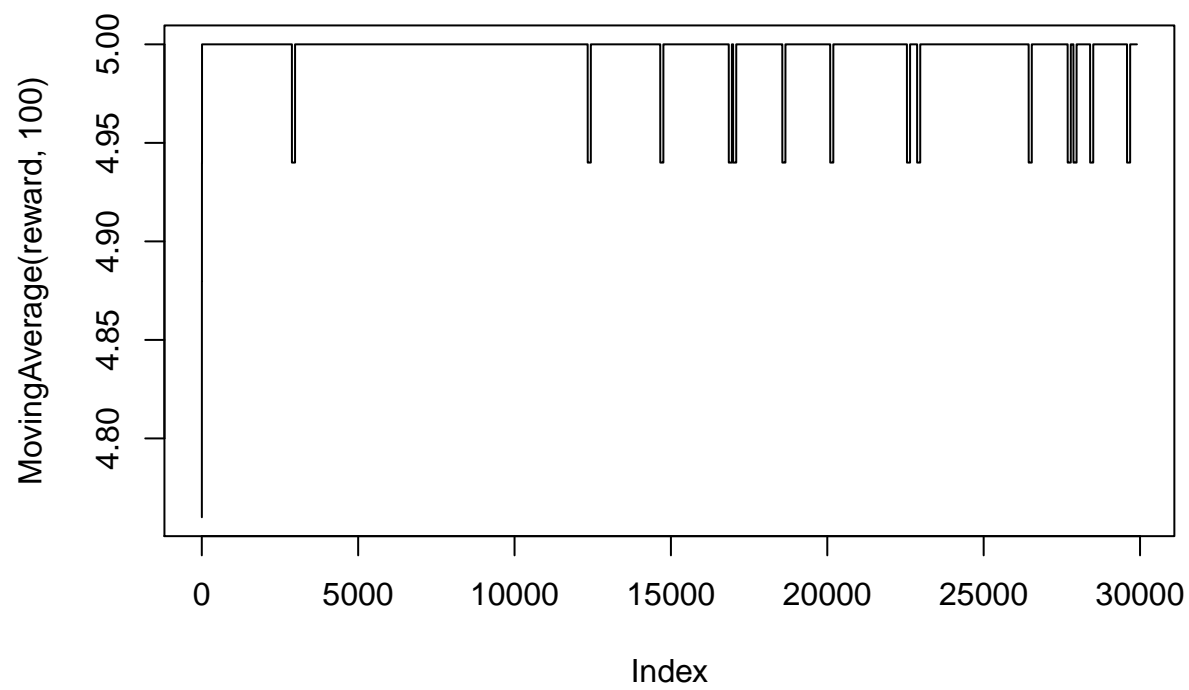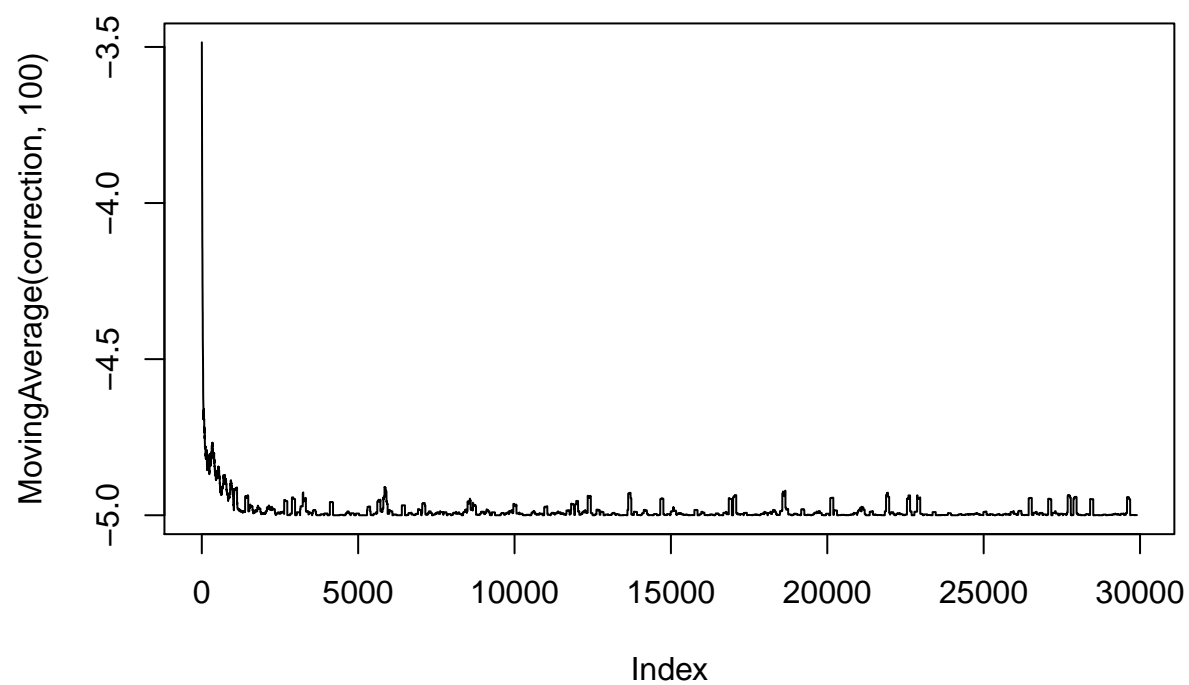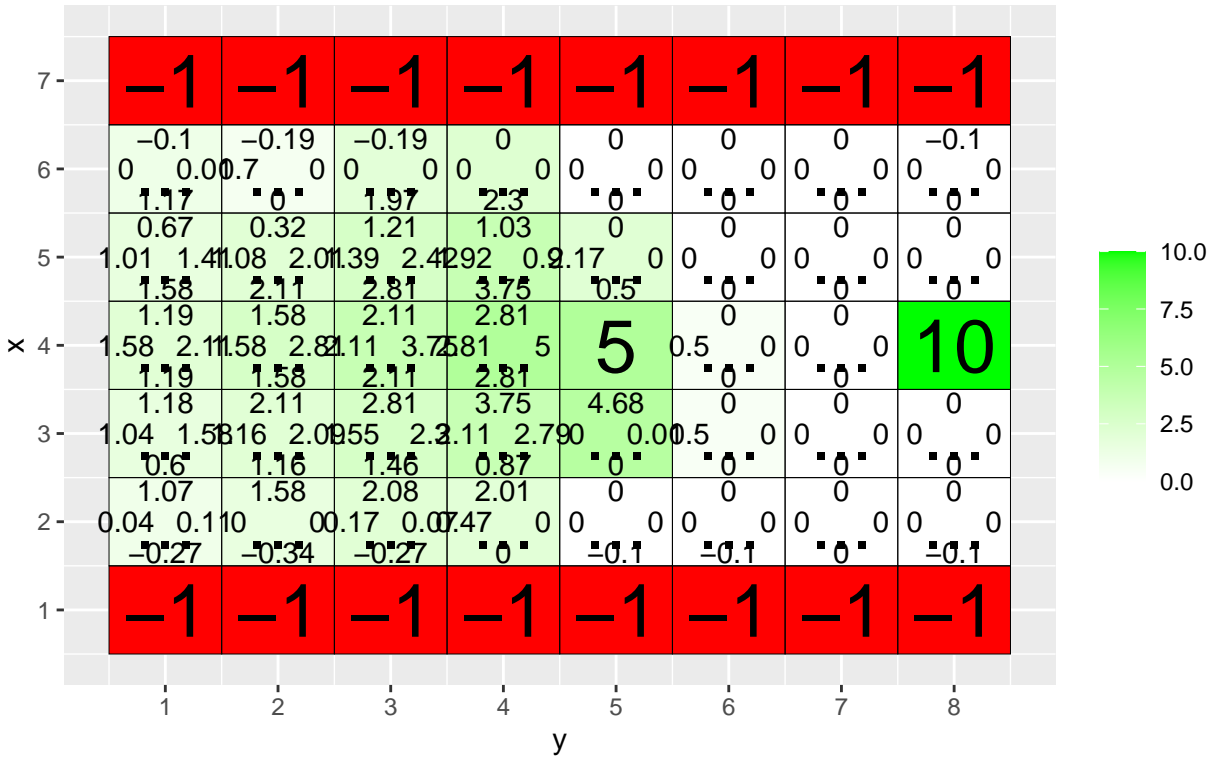
Q−table after 30000 iterations
(epsilon = 0.1 , alpha = 0.1 gamma = 0.5 , beta = 0 )

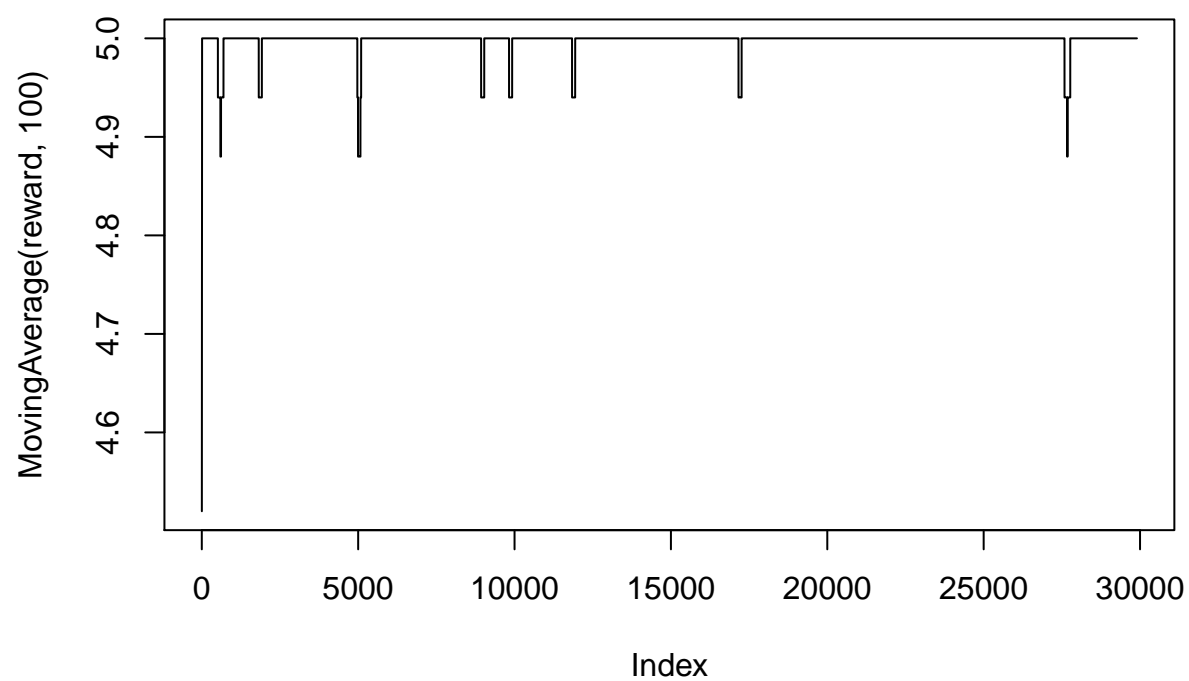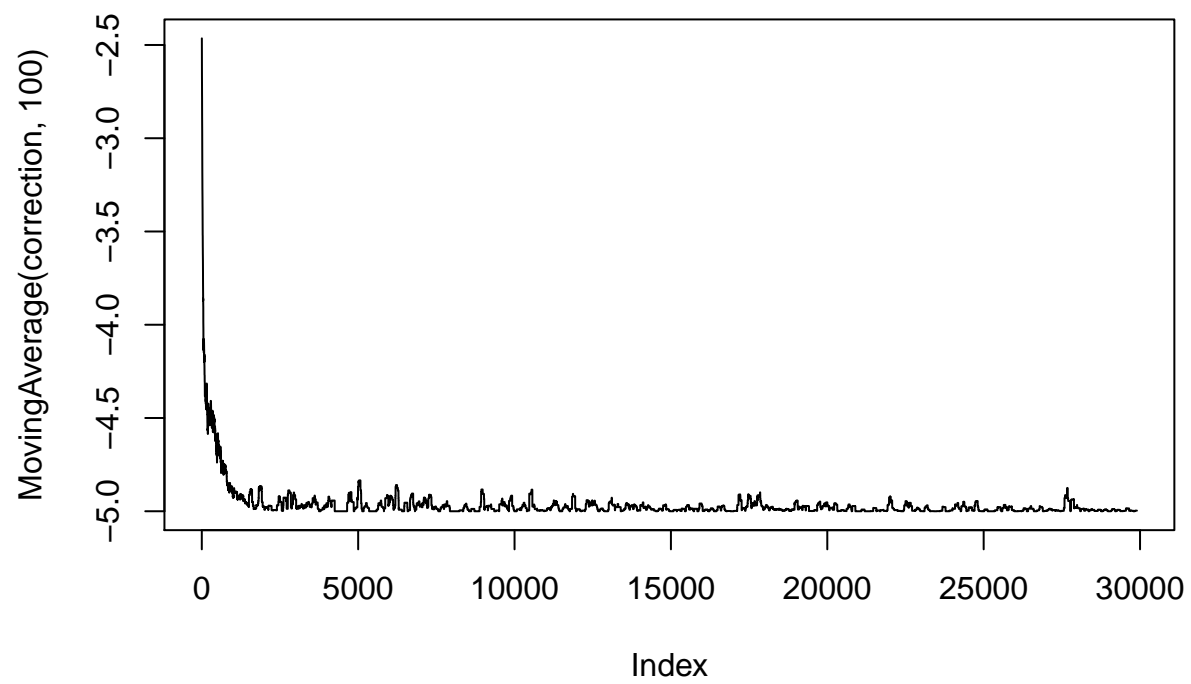| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 7 | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |
| 6 | −0.19 0 0 | −0.27 0 0.22 | −0.27 0.01 0.10 | −0.19 1.03 0 | 0 0 0 | −0.1 0 0 | −0.1 0 0 | −0.1 0 0 |
|  | 0.15 0 | 0 | 0.61 | 1.24 | 0 | 0 | 0 | 0 |
| 5 | 0.06 0.14 0.27 | 0.06 0.14 0.59 | 0.26 0.29 1.25 | 0.6 0.62 2.19 | 0 0.12 0 | 0 0 0 | 0 0 0 | 0 0 0 |
|  | 0.31 | 0.62 | 1.17 | 2.5 | 4.85 | 0.02 0 | 0 | 0 |
| 4 | 0.16 0.31 0.62 | 0.31 0.31 1.25 | 0.62 0.62 2.5 | 1.25 1.25 5 | 5 | 0 0.95 0 | 0 0 0 | 10 |
|  | 0.16 | 0.31 | 0.62 | 1.25 | 0.5 | 0 | 0 | 0 |
| 3 | 0.31 0.13 0.26 | 0.62 0.15 0.56 | 1.25 0.26 1.10 | 2.5 1.58 0.41 | 0.5 1.11 0 | 0 0 | 0 0 | 0 0 |
|  | 0.06 | 0.14 | 0.26 | 0.34 | 0 | 0 | 0 | 0 |
| 2 | 0.03 0.01 0.14 0 | 0.31 0 | 0.6 0.02 0.05 | 1.04 0.07 0 | 0 0 0 | 0 0 0 | 0 0 0 | 0 0 0 |
|  | −0.1 | −0.19 | −0.1 | −0.1 | 0 | 0 | 0 | 0 |
| 1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |

x

y

10.0
7.5
5.0
2.5
0.0

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

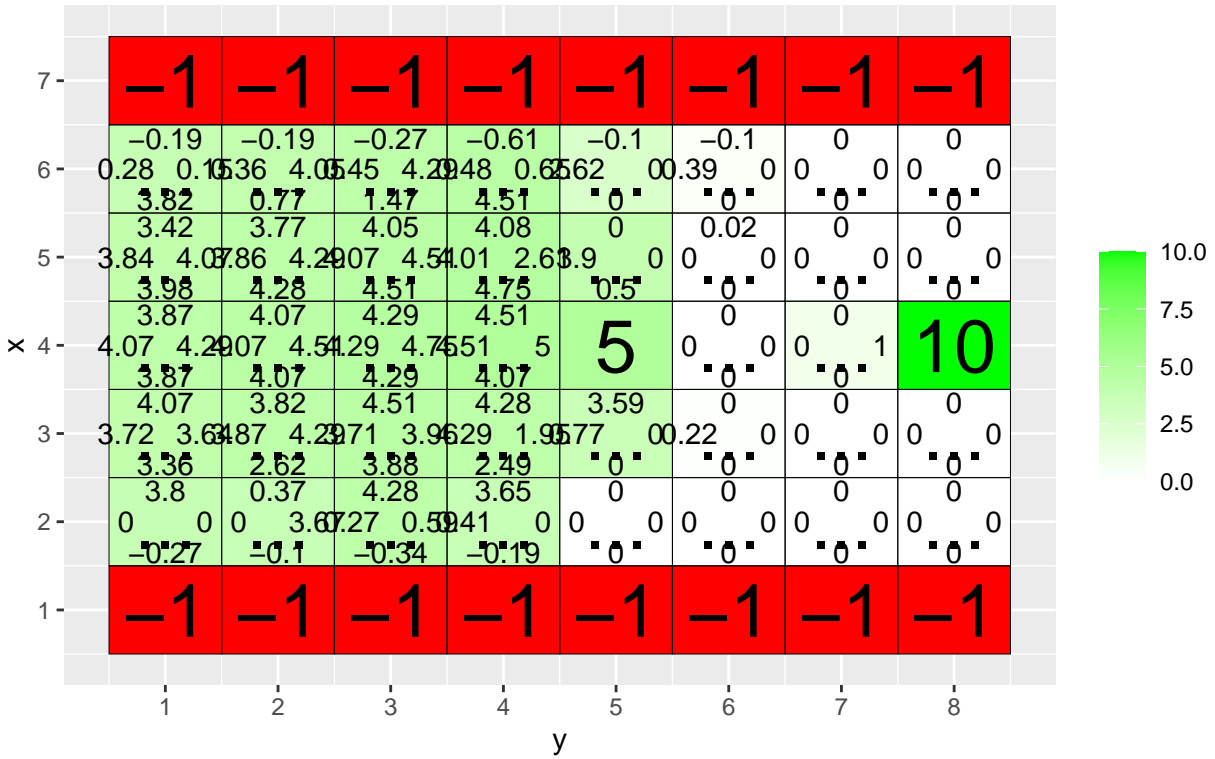| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 7 | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |
| 6 | −0.1<br>0   0.0<br>1.17 | −0.19<br>0.7   0<br>0 | −0.19<br>0   0<br>1.97 | 0<br>0   0<br>2.3 | 0<br>0   0<br>0 | 0<br>0   0<br>0 | 0<br>0   0<br>0 | −0.1<br>0   0<br>0 |
| 5 | 0.67<br>1.01  1.41<br>1.58 | 0.32<br>1.08  2.01<br>2.11 | 1.21<br>1.39  2.42<br>2.81 | 1.03<br>2.92  0.2<br>3.75 | 0<br>2.17   0<br>0.5 | 0<br>0   0<br>0 | 0<br>0   0<br>0 | 0<br>0   0<br>0 |
| 4 | 1.19<br>1.58  2.11<br>1.19 | 1.58<br>1.58  2.81<br>1.58 | 2.11<br>2.11  3.75<br>2.11 | 2.81<br>2.81  5<br>2.81 | 5 | 0<br>0.5   0<br>0 | 0<br>0   0<br>0 | 10 |
| 3 | 1.18<br>1.04  1.58<br>0.6 | 2.11<br>2.16  2.09<br>1.16 | 2.81<br>1.55  2.3<br>1.46 | 3.75<br>3.11  2.79<br>0.87 | 4.68<br>0   0.0<br>0 | 0<br>0.5   0<br>0 | 0<br>0   0<br>0 | 0<br>0   0<br>0 |
| 2 | 1.07<br>0.04  0.11<br>−0.27 | 1.58<br>10   0<br>−0.34 | 2.08<br>0.17  0.0<br>−0.27 | 2.01<br>7.47   0<br>0 | 0<br>0   0<br>−0.1 | 0<br>0   0<br>−0.1 | 0<br>0   0<br>0 | 0<br>0   0<br>−0.1 |
| 1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |

x

y

Legend: 10.0, 7.5, 5.0, 2.5, 0.0

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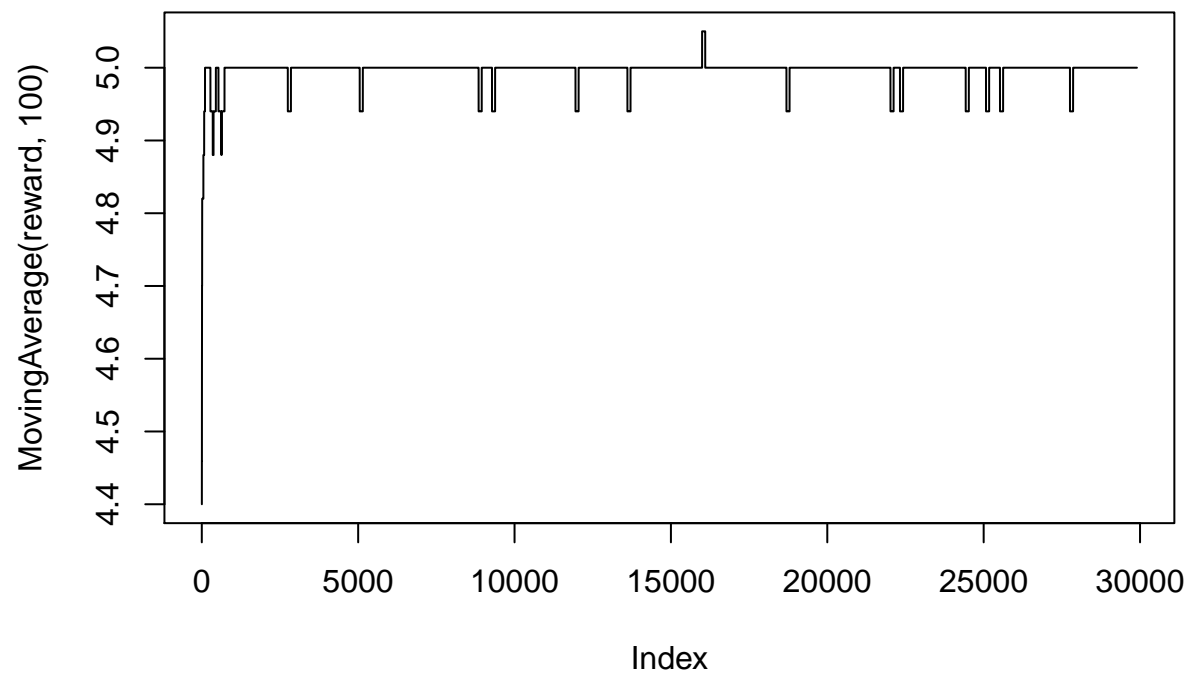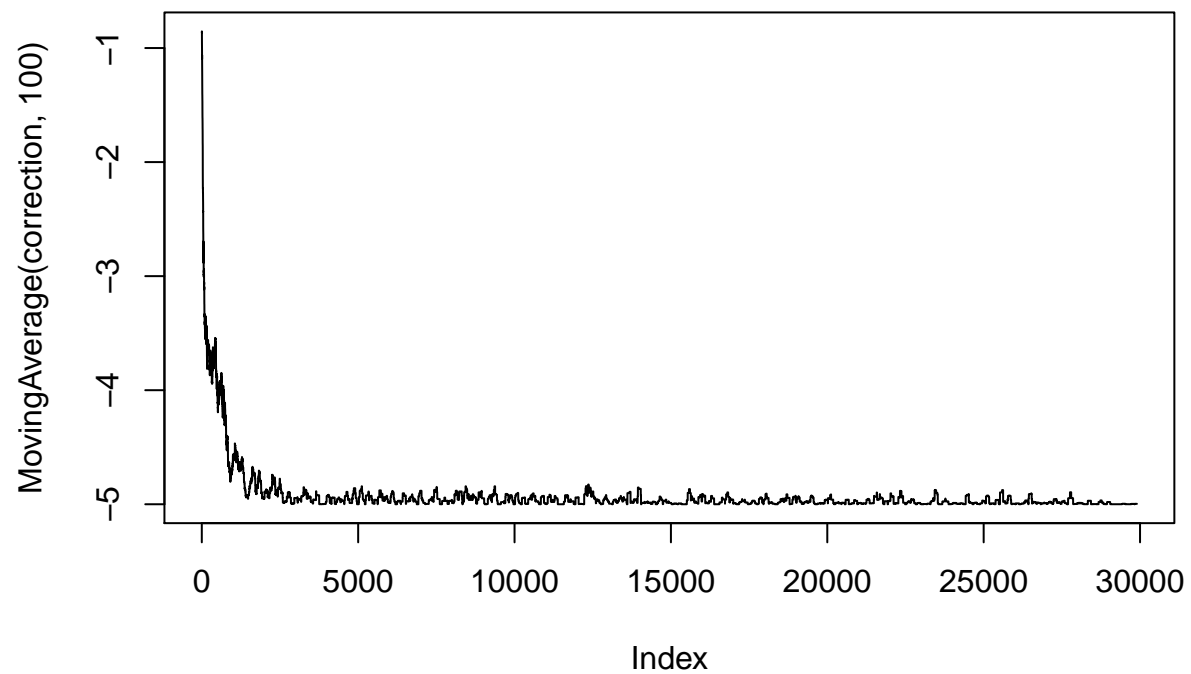## (epsilon =  0.1 , alpha =  0.1 gamma =  0.95 , beta =  0 )

x

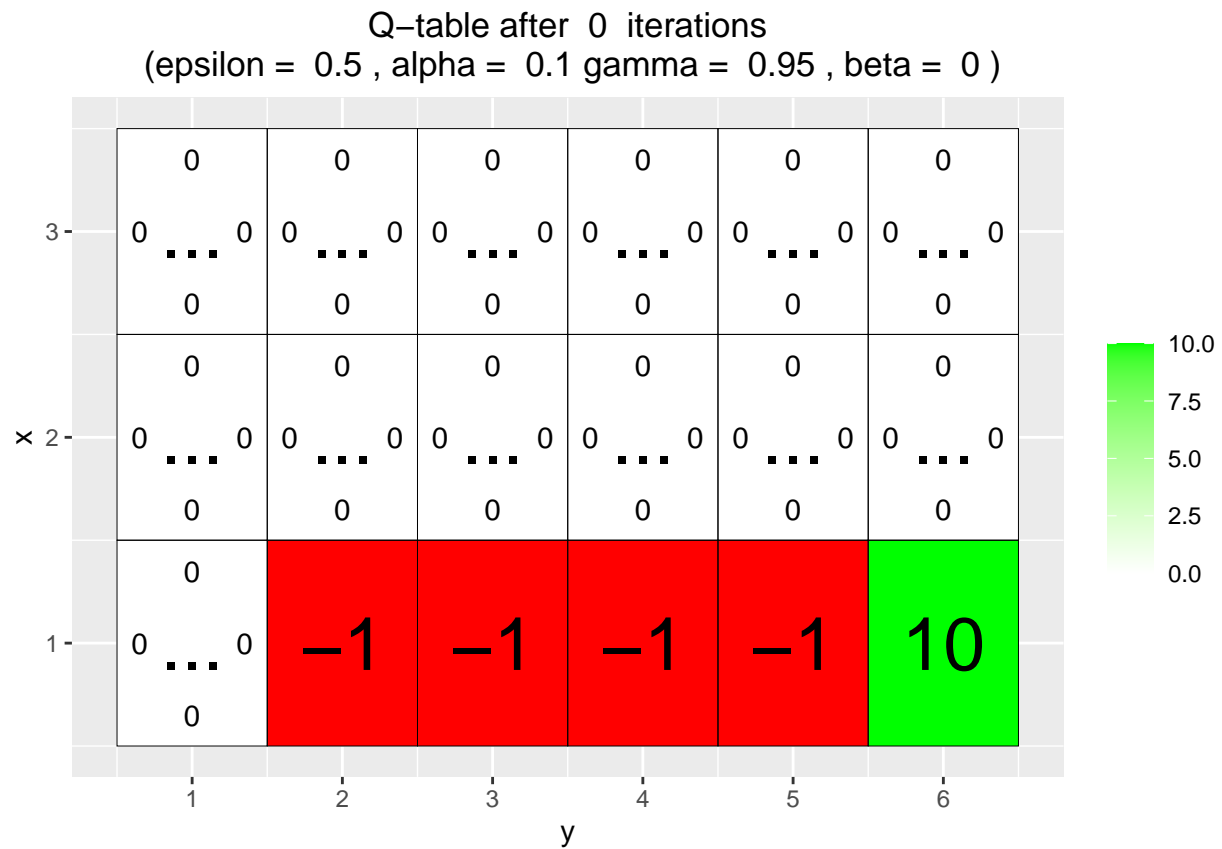| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 7 | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |
| 6 | −0.19 / 0.28  0.10 / 3.82 | −0.19 / 3.36  4.05 / 0.77 | −0.27 / 3.45  4.29 / 1.47 | −0.61 / 3.48  0.65 / 4.51 | −0.1 / 3.62   0 / 0 | −0.1 / 0.39   0 / 0 | 0 / 0   0 / 0 | 0 / 0   0 / 0 |
| 5 | 3.42 / 3.84  4.07 / 3.98 | 3.77 / 3.86  4.29 / 4.28 | 4.05 / 4.07  4.54 / 4.51 | 4.08 / 4.01  2.63 / 4.75 | 0 / 3.9   0 / 0.5 | 0.02 / 0   0 / 0 | 0 / 0   0 / 0 | 0 / 0   0 / 0 |
| 4 | 3.87 / 4.07  4.29 / 3.87 | 4.07 / 4.07  4.54 / 4.07 | 4.29 / 4.29  4.74 / 4.29 | 4.51 / 4.51  5 / 4.07 | 5 | 0 / 0   0 / 0 | 0 / 0   1 / 0 | 10 |
| 3 | 4.07 / 3.72  3.64 / 3.36 | 3.82 / 4.87  4.23 / 2.62 | 4.51 / 3.71  3.96 / 3.88 | 4.28 / 3.29  1.95 / 2.49 | 3.59 / 3.77   0 / 0 | 0 / 0.22   0 / 0 | 0 / 0   0 / 0 | 0 / 0   0 / 0 |
| 2 | 3.8 / 0   0 / −0.27 | 0.37 / 0  3.67 / −0.1 | 4.28 / 7.27  0.59 / −0.34 | 3.65 / 3.41   0 / −0.19 | 0 / 0   0 / 0 | 0 / 0   0 / 0 | 0 / 0   0 / 0 | 0 / 0   0 / 0 |
| 1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |

y

Legend: 10.0 / 7.5 / 5.0 / 2.5 / 0.0

```
# Environment C (the effect of beta).

H <- 3
W <- 6

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

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

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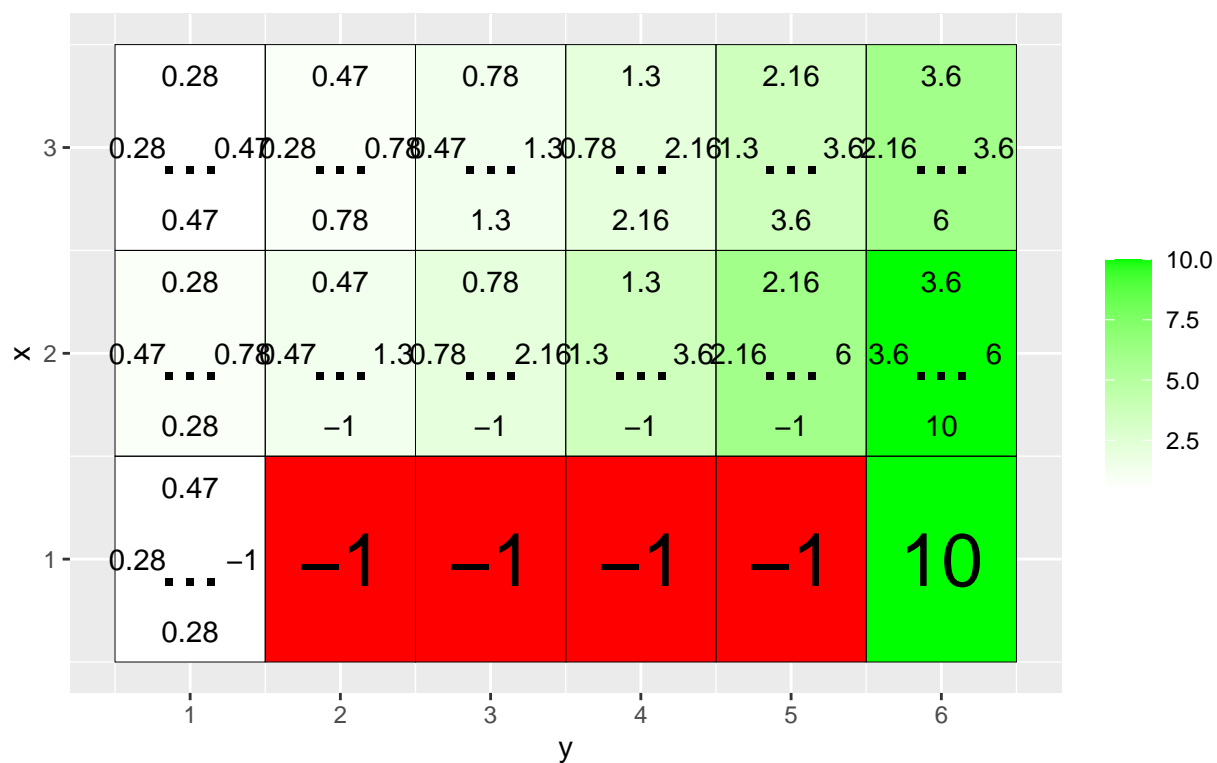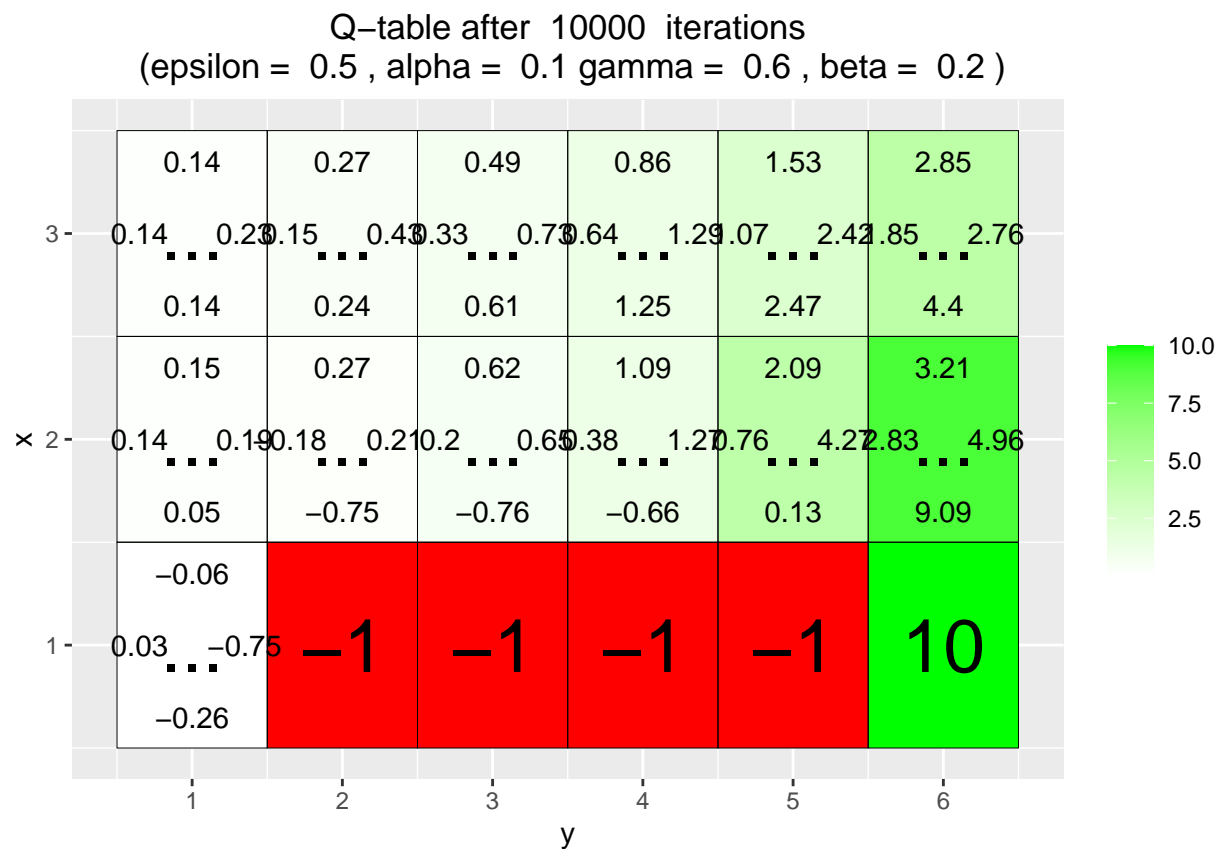
## Q–table after  0  iterations
## (epsilon =  0.5 , alpha =  0.1 gamma =  0.95 , beta =  0 )



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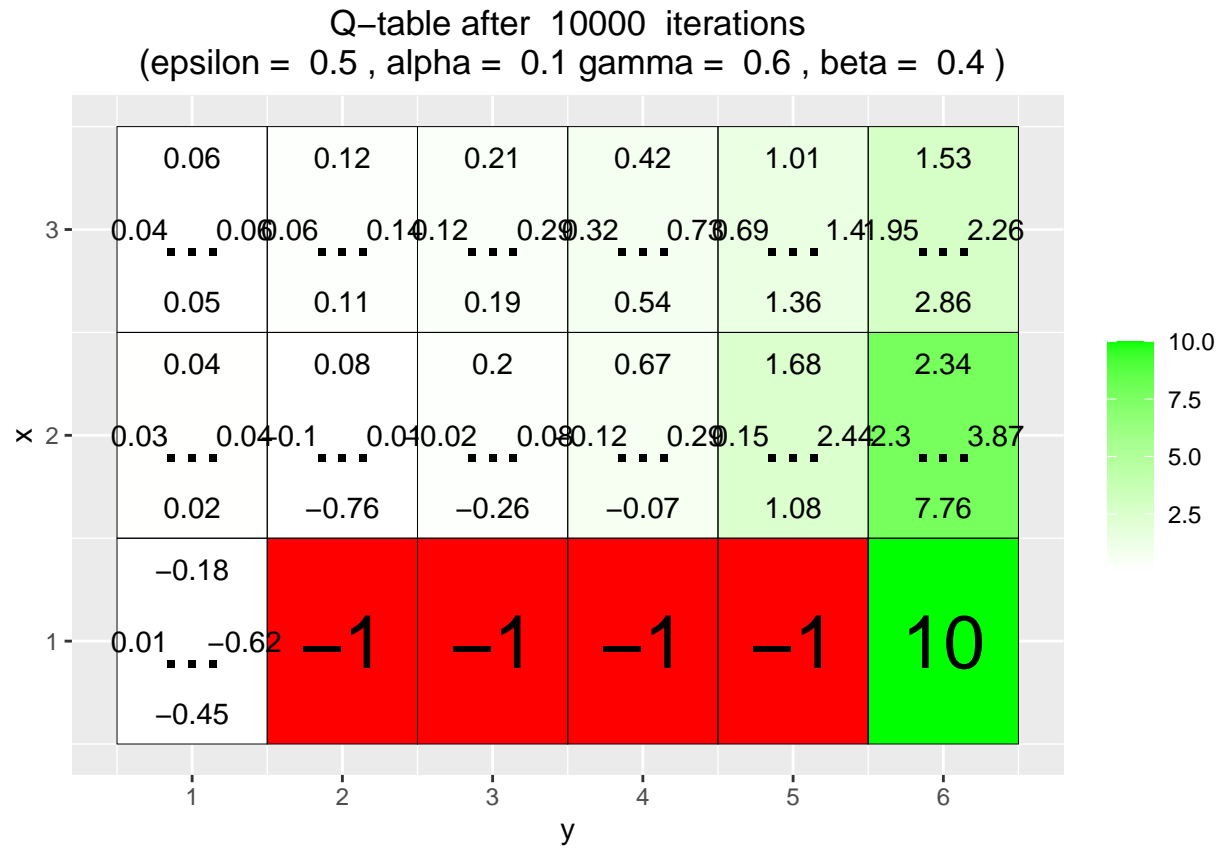
# Q−table after  10000  iterations
## (epsilon =  0.5 , alpha =  0.1 gamma =  0.6 , beta =  0 )

## 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 )