# Lab 3 RL Individual Report

Akshay Gurudath

29/09/2021

## 2.1 Q-Learning Algorithm

```r
# By Jose M. PeÃ±a and Joel Oskarsson.
# For teaching purposes.
# jose.m.pena@liu.se.


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

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

```
## Warning: package 'ggplot2' was built under R version 4.0.5
```

```r
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){

  # 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))
```

```r
  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) +
          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.
  max <- which(q_table[x,y,] == max(q_table[x,y,]))
  if(length(max) > 1){
    max <- sample(max, 1)
  }
  return(max)

}

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

  # Get an epsilon-greedy action for state (x,y) from q_table.
  #
  # Args:
```

2

```r
  #   x, y: state coordinates.
  #   epsilon: probability of acting randomly.
  #
  # Returns:
  #   An action, i.e. integer in {1,2,3,4}.

  # Your code here.
  a=GreedyPolicy(x,y)

  all=c(1,2,3,4)

  if(runif(1)<=1-epsilon){
    return(a)
  }
  else{
    return(sample(all,1))
  }

}

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:
```

```r
#    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 <<-.

# Your code here.
#start_state=c(sample.int(5,1),sample.int(7,1))
x=start_state[1]
y=start_state[2]

episode_correction=0

repeat{
  # Follow policy, execute action, get reward.
  a=EpsilonGreedyPolicy(x,y,epsilon=epsilon)
  #a_greedy=GreedyPolicy(x,y)

  new_state=transition_model(x,y,action=a,beta=beta)
 # print(new_state)
  reward=reward_map[new_state[1],new_state[2]]
  temp_diff=-q_table[x,y,a]+(reward+gamma*max(q_table[new_state[1],new_state[2],]))

  q_table[x,y,a]<<-q_table[x,y,a]+alpha*temp_diff
  # Q-table update.
  x=new_state[1]
  y=new_state[2]
  episode_correction=episode_correction+temp_diff
  if(reward!=0)
    # End episode.
    return (c(reward,episode_correction))
  }

}
```

## 2.2 Environment A

```r
###############################################################################################
# Q-Learning Environments
###############################################################################################

# Environment A (learning)
set.seed(12345)
H <- 5
```
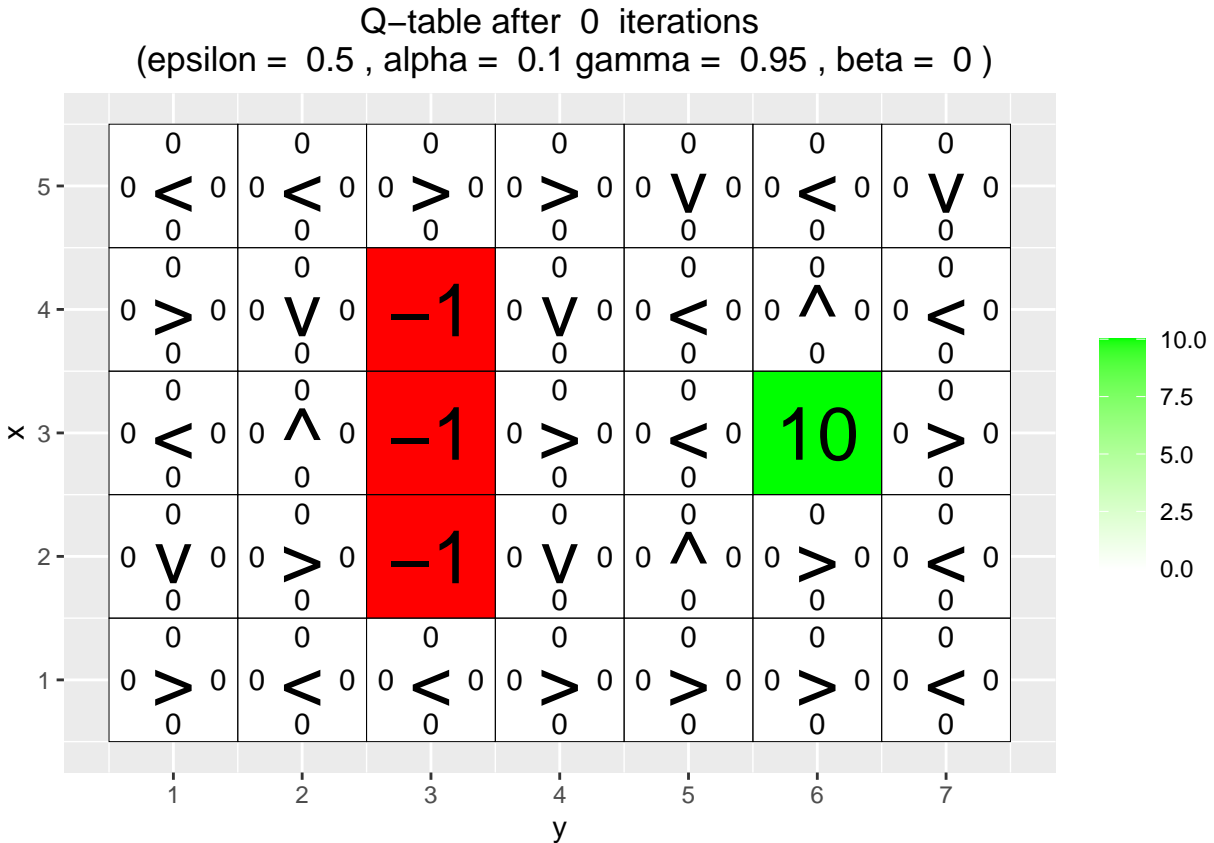
4

```
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()
```
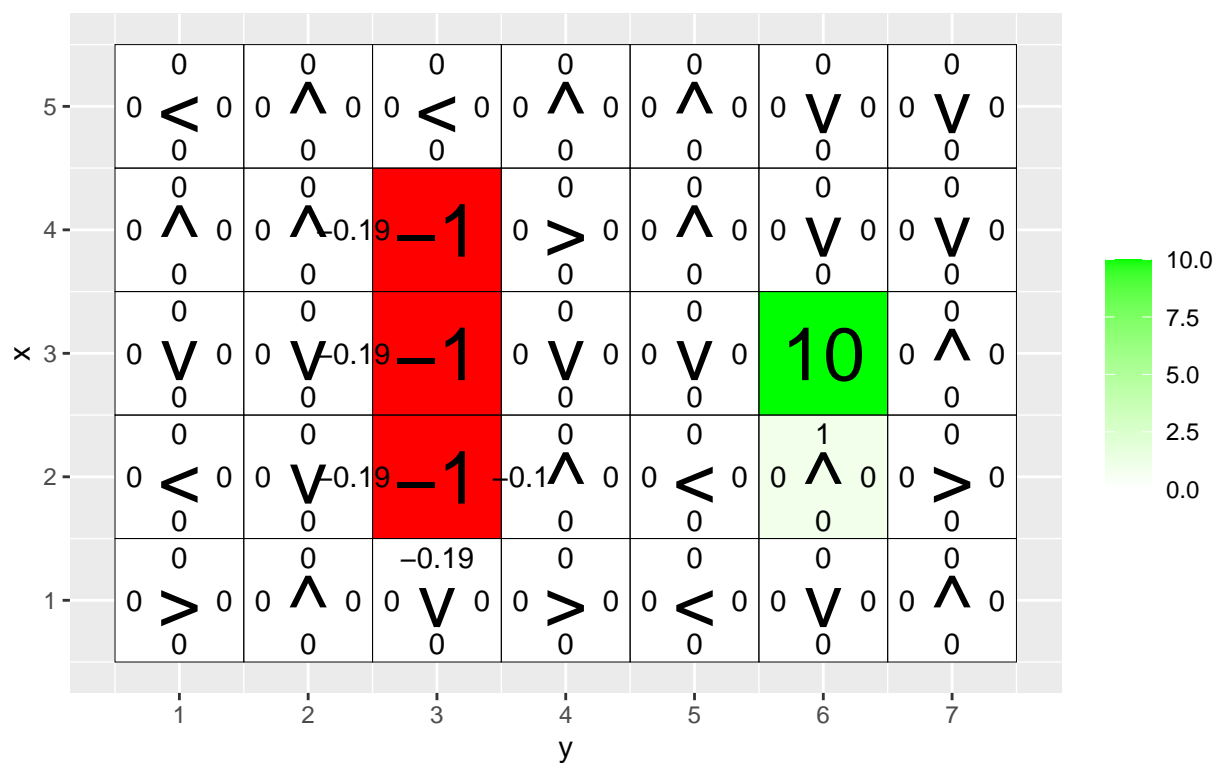


Q–table after 0 iterations
(epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0 )

```
for(i in 1:10000){
  foo <- q_learning(start_state = c(3,1),epsilon=0.5,gamma=0.95)

  if(any(i==c(10,100,1000,10000,20000,30000,40000)))
    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 )

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

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

a. The agent has learned almost nothing in the first 10 episodes. It has learned that the states close to -1 are states to stay away from. However, because of the duration, it has only learned so much and not explored other paths.

b. The final greedy policy isn't exactly optimal for all states as some states are less visited than the others and their q-values cannot be trusted with the same certainty as some other states which are visited more often.

c. No multiple paths are not reflected in the current solution (both upward and downward). These paths are not perfect and a better way to show multiple paths is by initializing the states randomly each time an episode starts. This will result in a more uniform visitation of states and thus will show other paths to the objective from every state.

## 2.3 Environment B

```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
```

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

vis_environment()
```
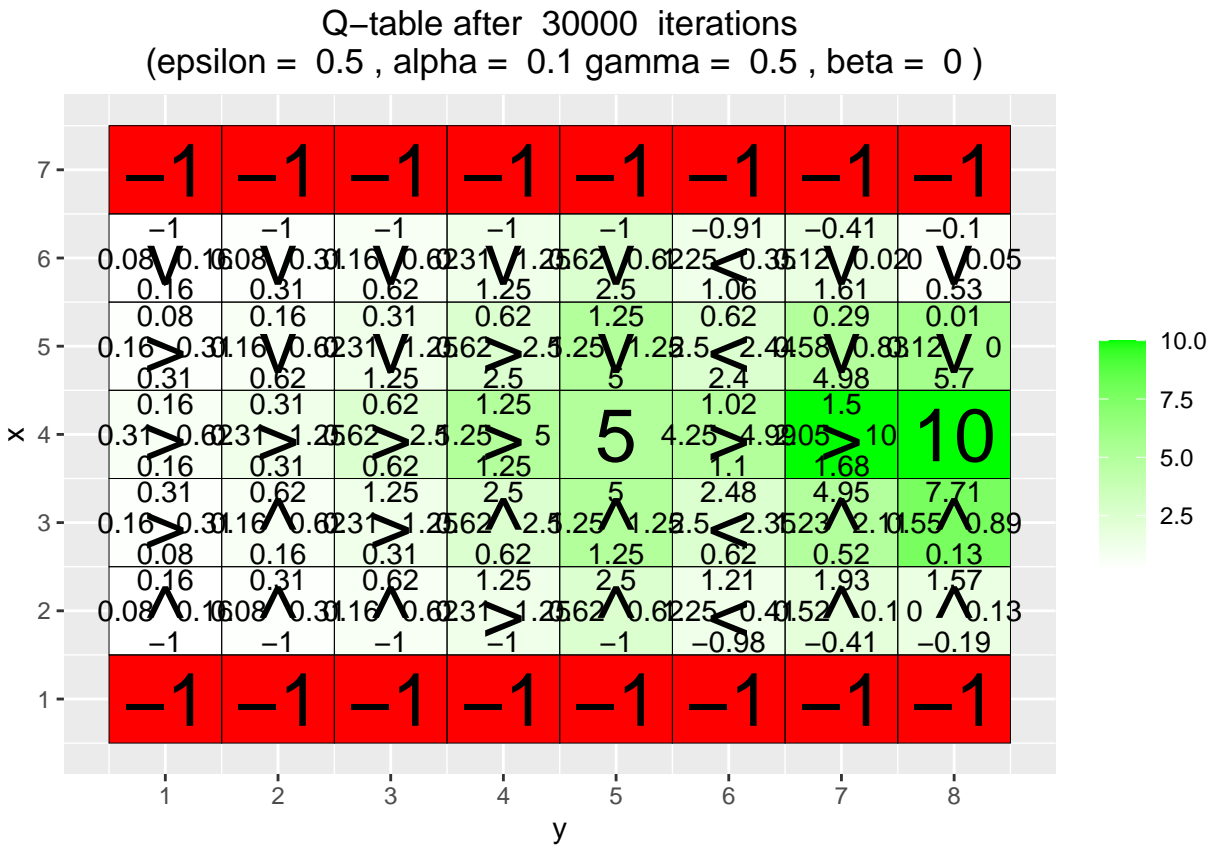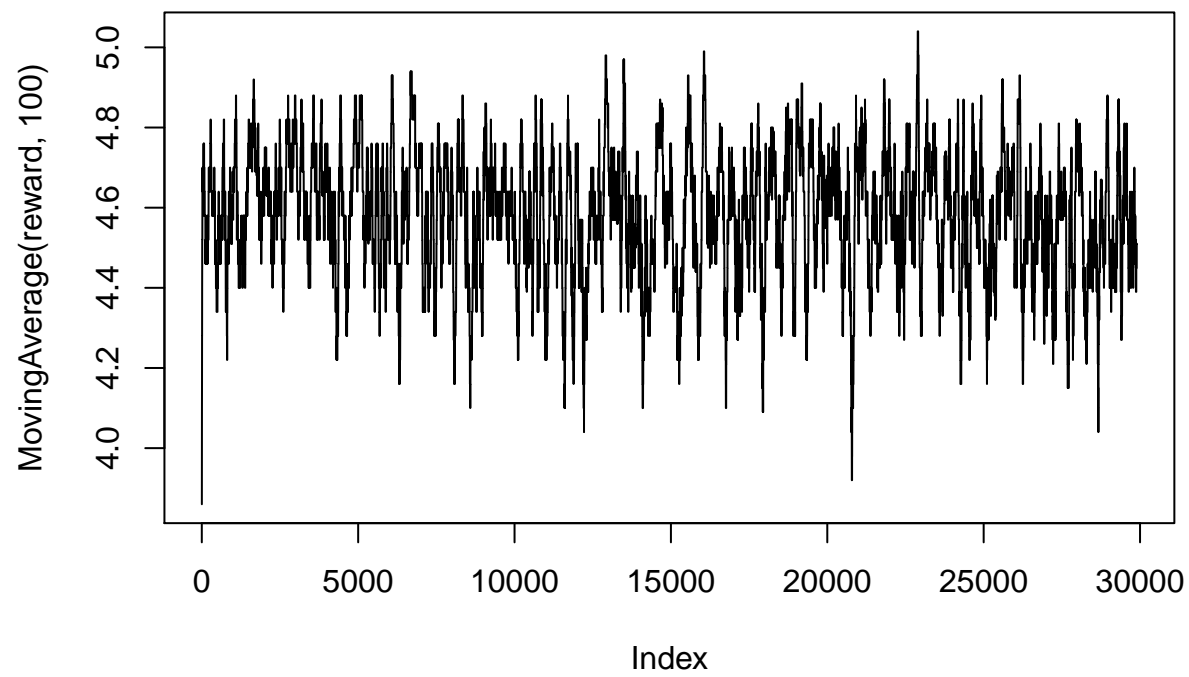


MovingAverage <- function(x, n){

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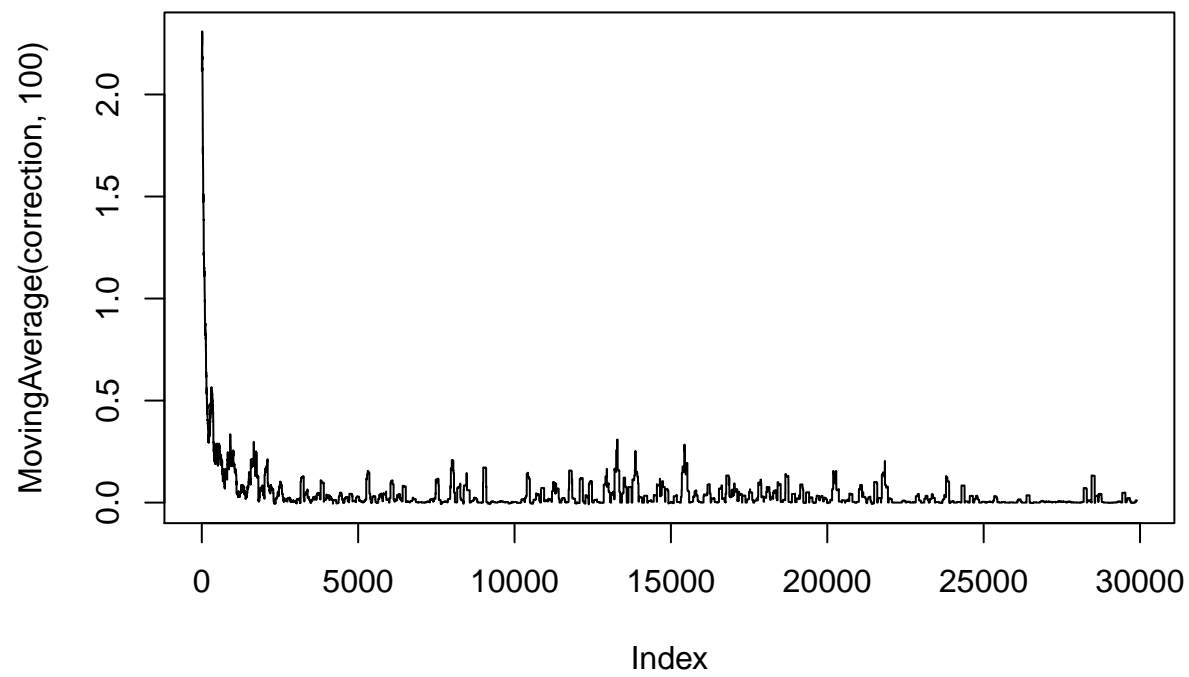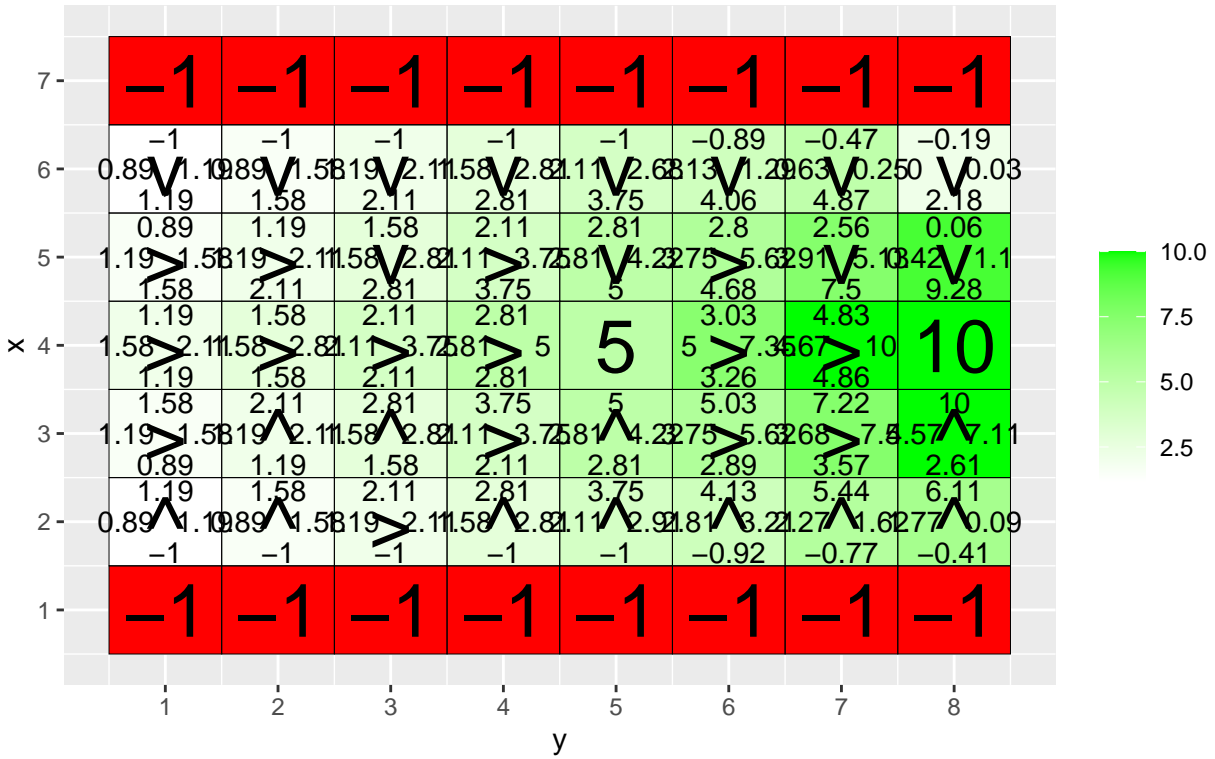
  return (rsum)
}

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")
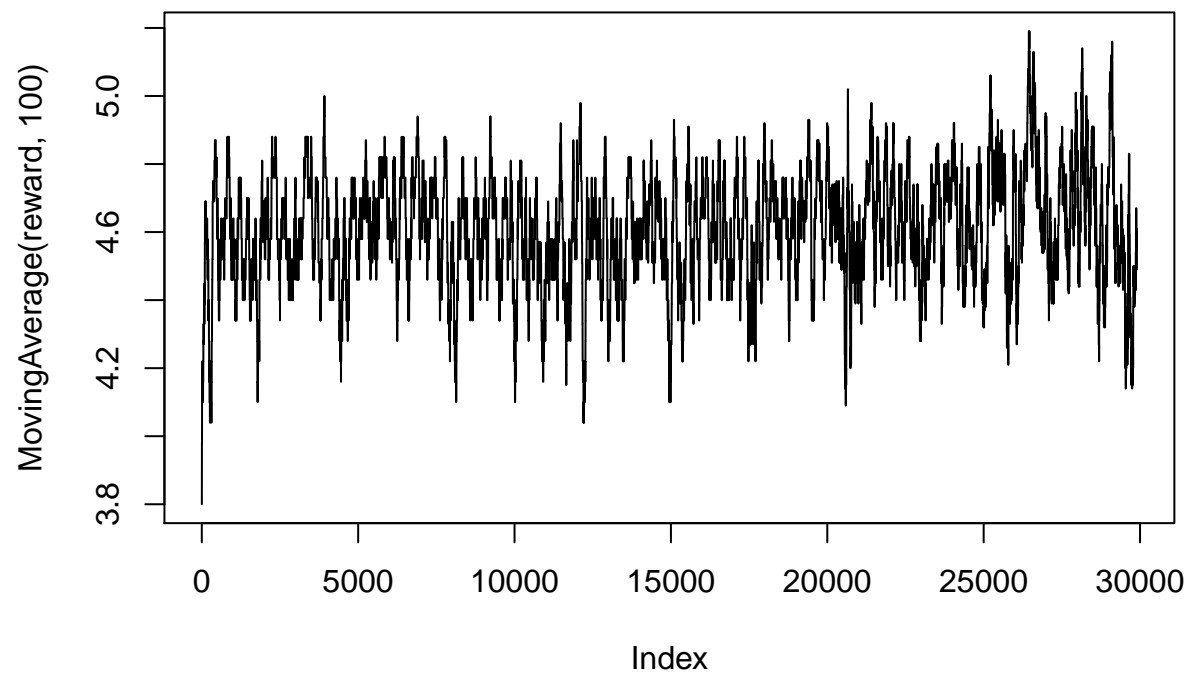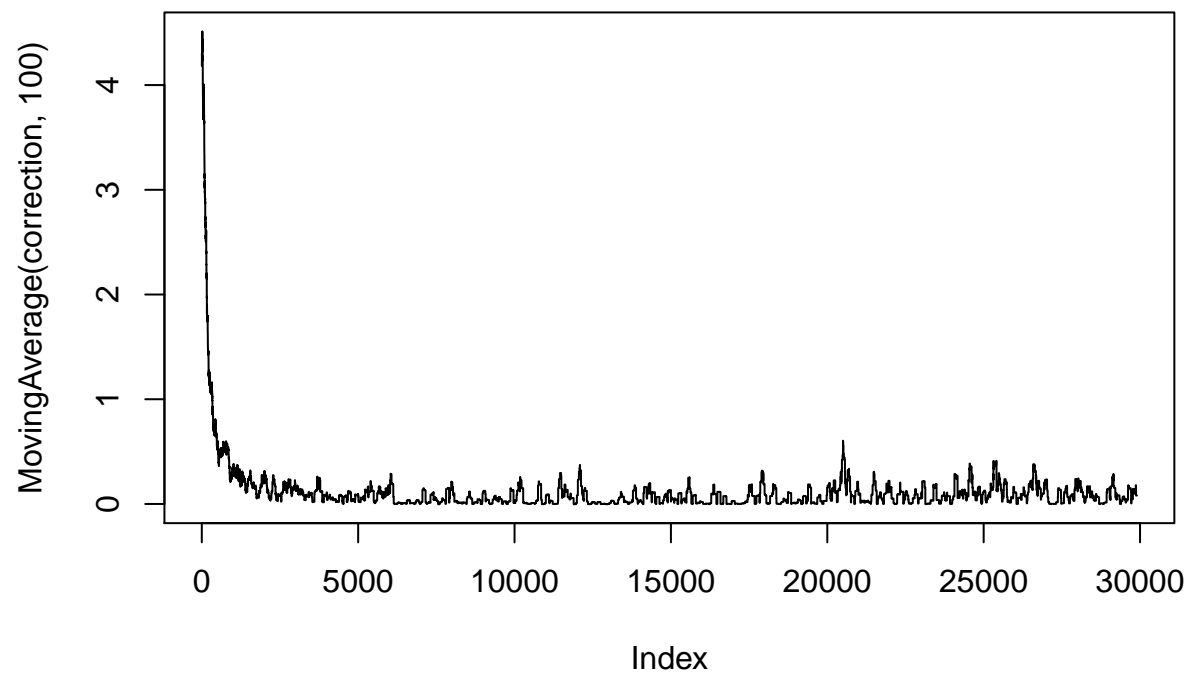  plot(MovingAverage(correction,100),type = "l")
```

```
}
```



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 )

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

When $\gamma = 0.5$, we notice that there are a few states close to 10 which point to 10, and many other states seem to point towards 5. We also notice that a lot of these states are not colored, which means that they don't have significant values of q. This is probably because states far away from the reward states get a low cumulative reward as the weighting factor is less (0.5).

When $\gamma = 0.75$, we notice that a lot more of these states are colored green, because they get a better cumulative reward. We also notice that there is more tendency to move to the +10 reward.

When $\gamma = 0.95$ most of the states are colored green implying that they get the best cumulative rewards because of the best discounting. We also notice that there is most tendency to go to the +10 reward. We also see a jump in the cumulative reward showing that the agent now is tending to go to the +10 reward states,

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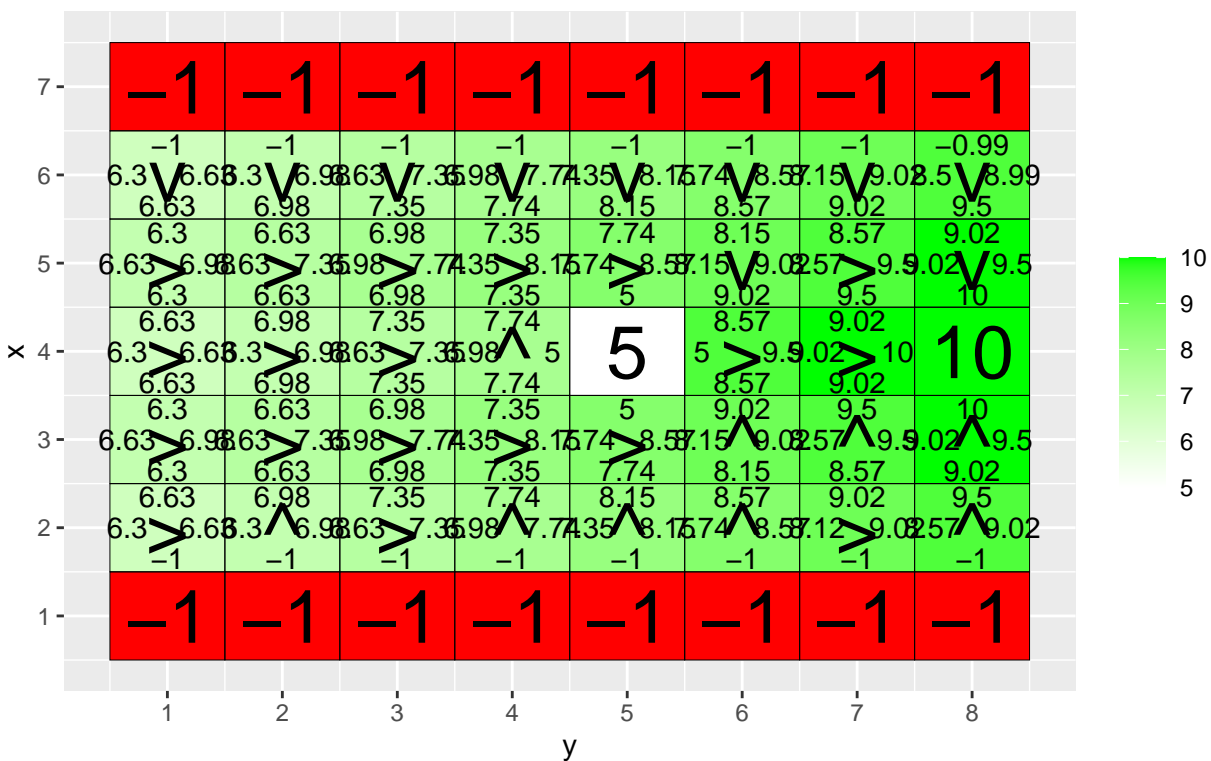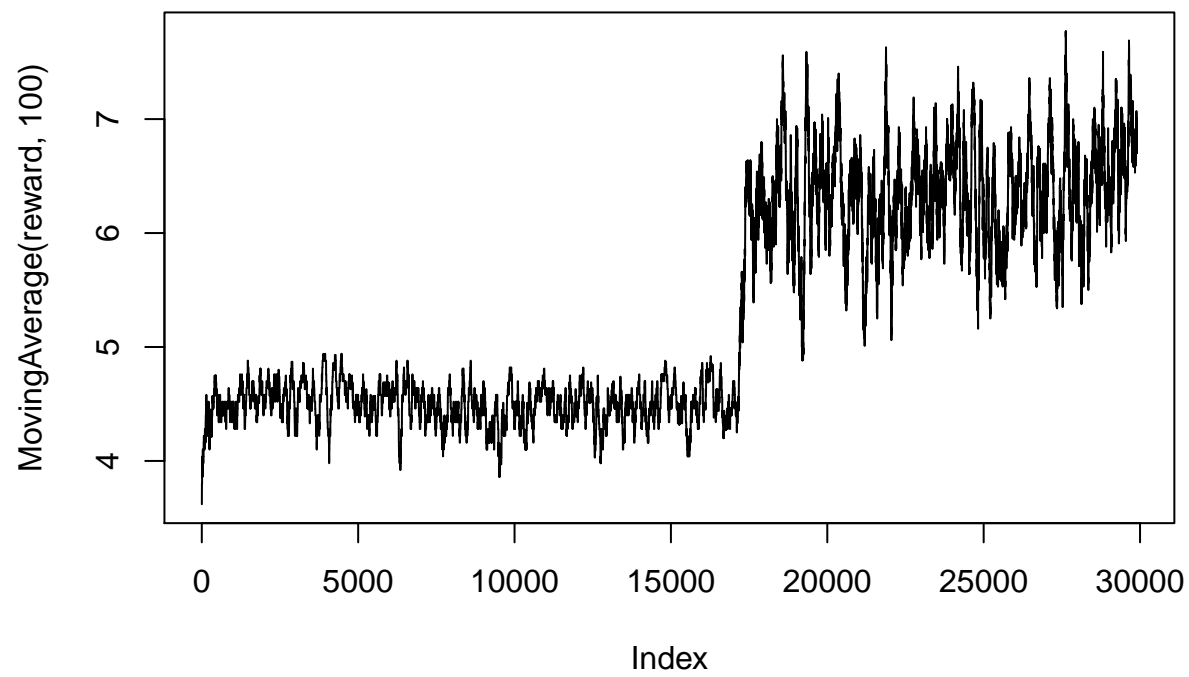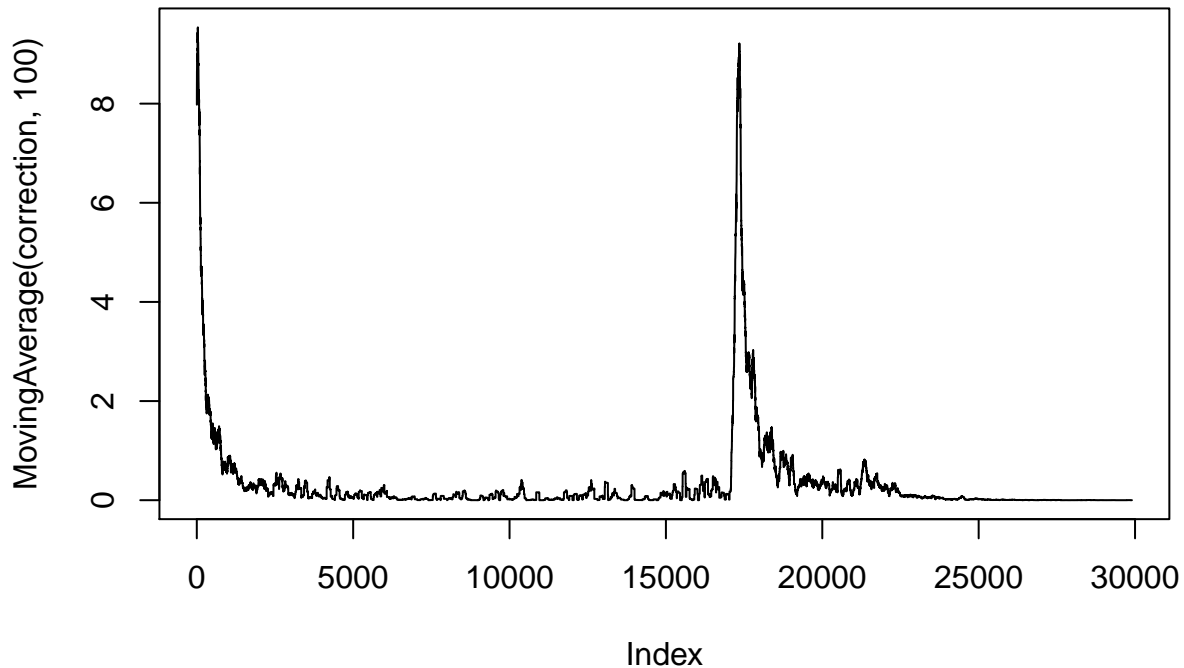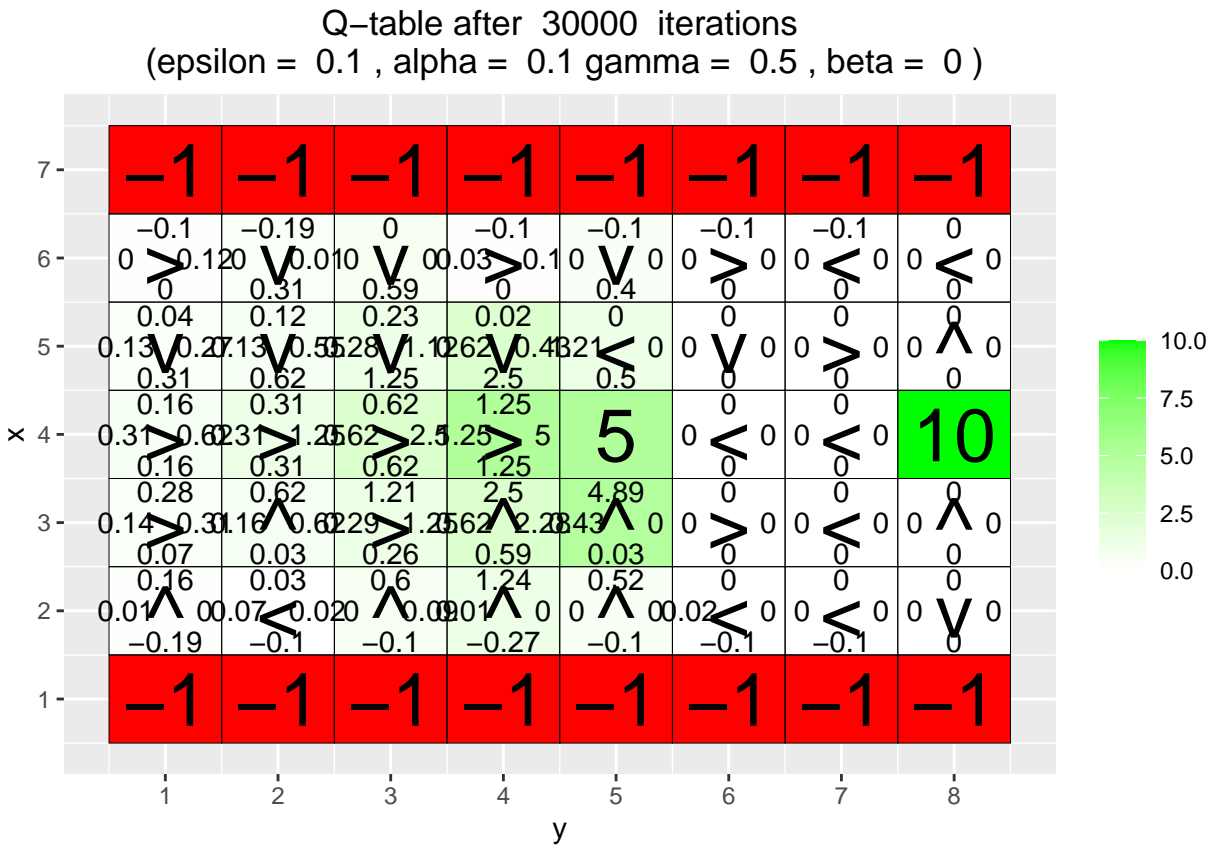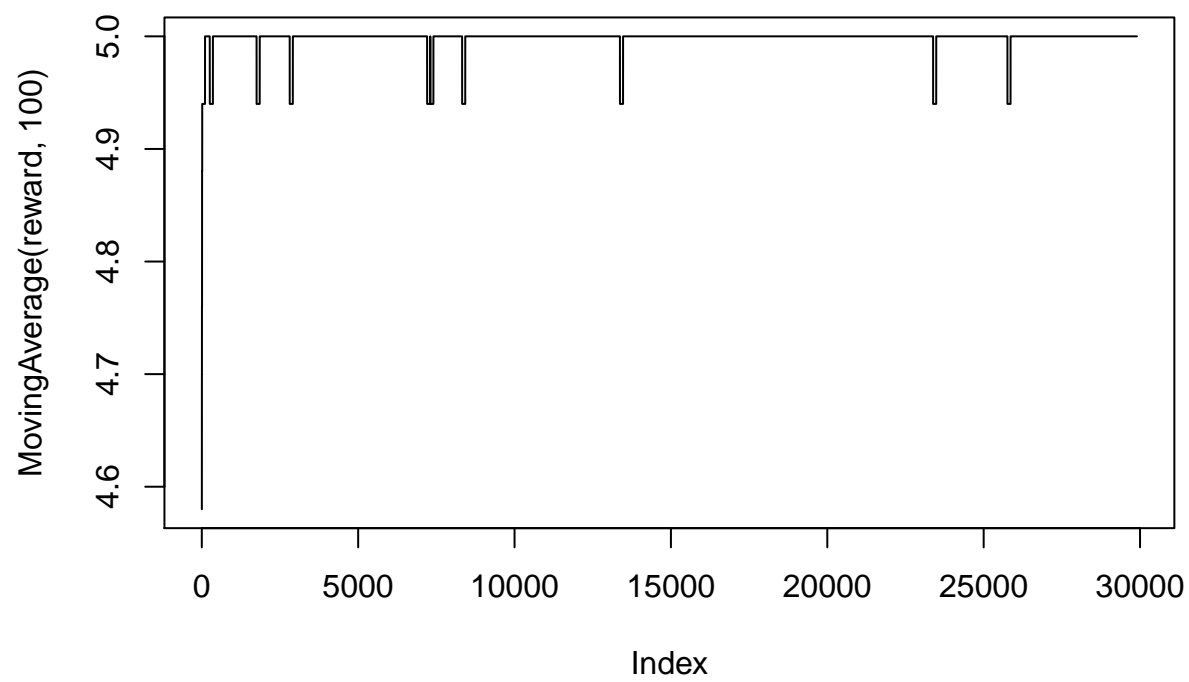for(j in c(0.5,0.75,0.95)){
  q_table <- array(0,dim = c(H,W,4))
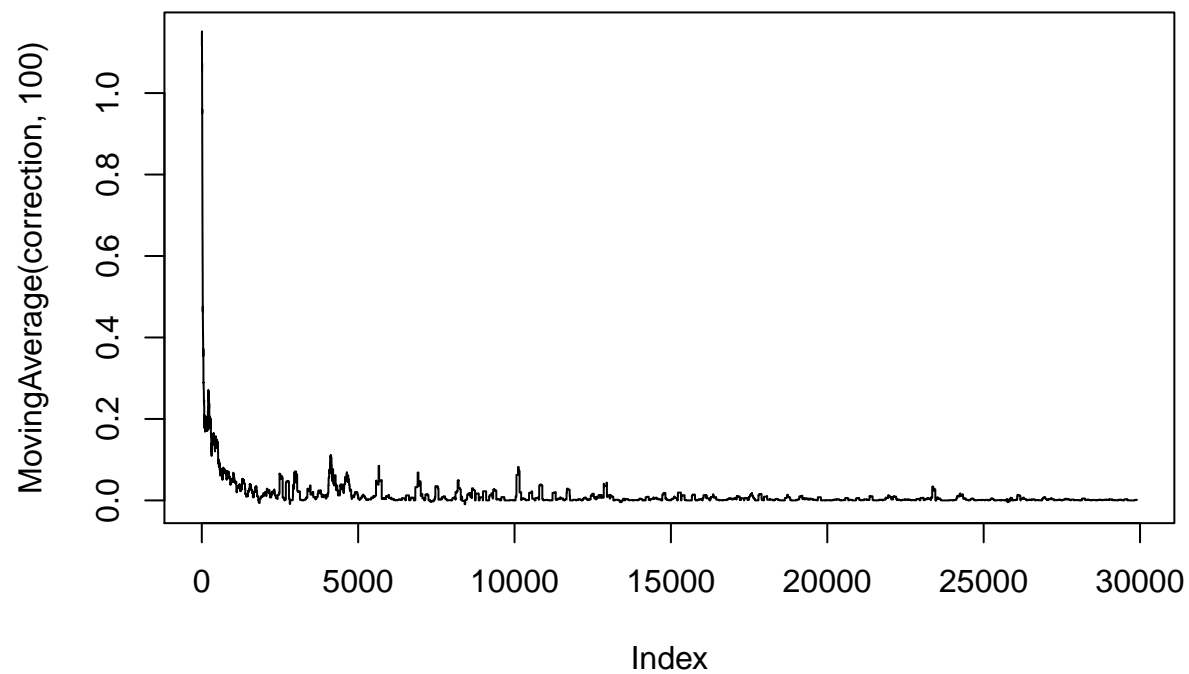  reward <- NULL
  correction <- NULL

  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")
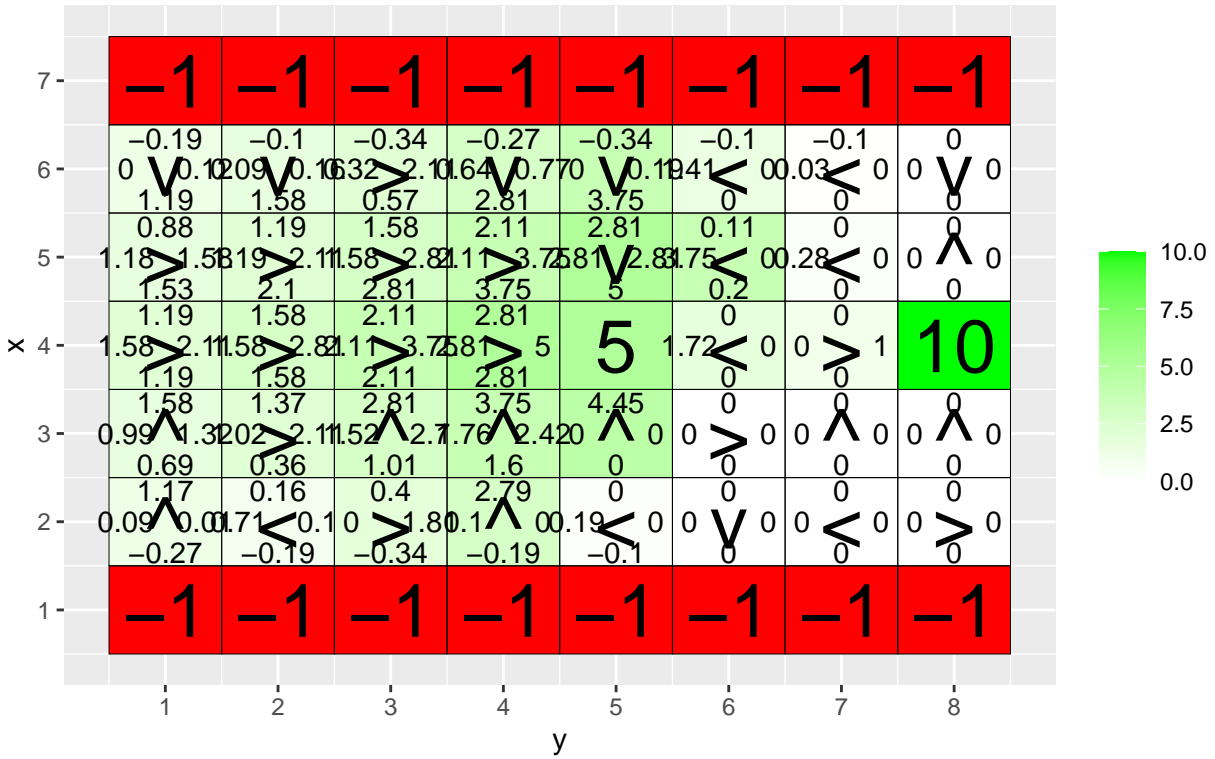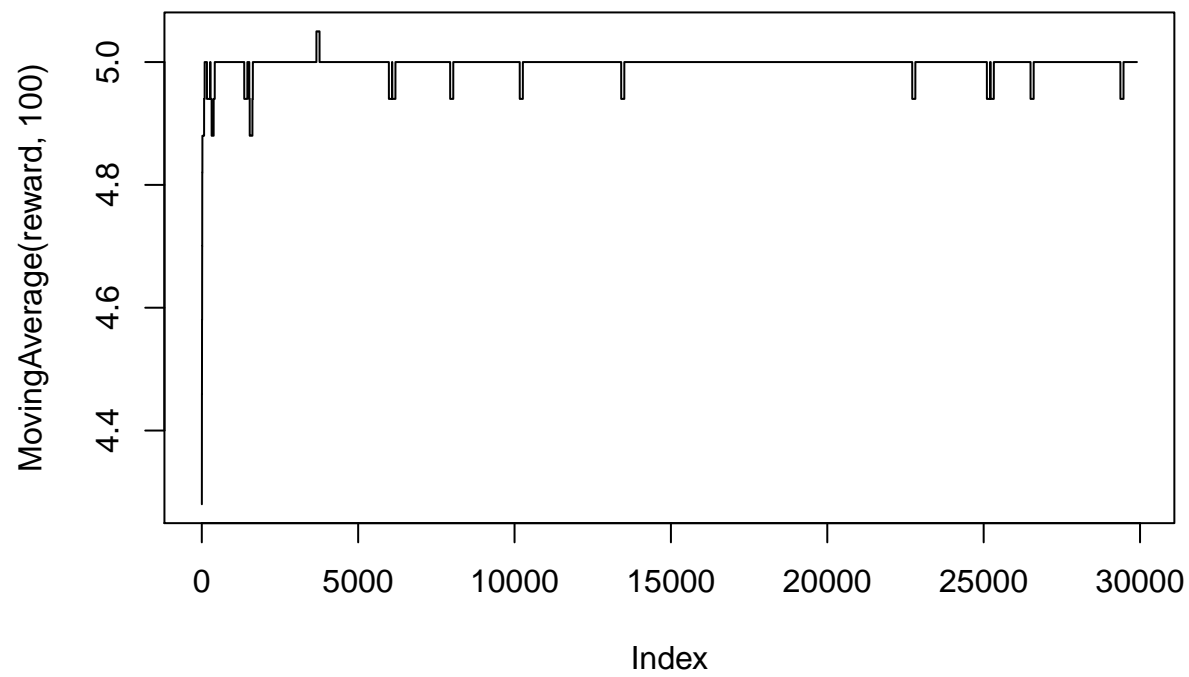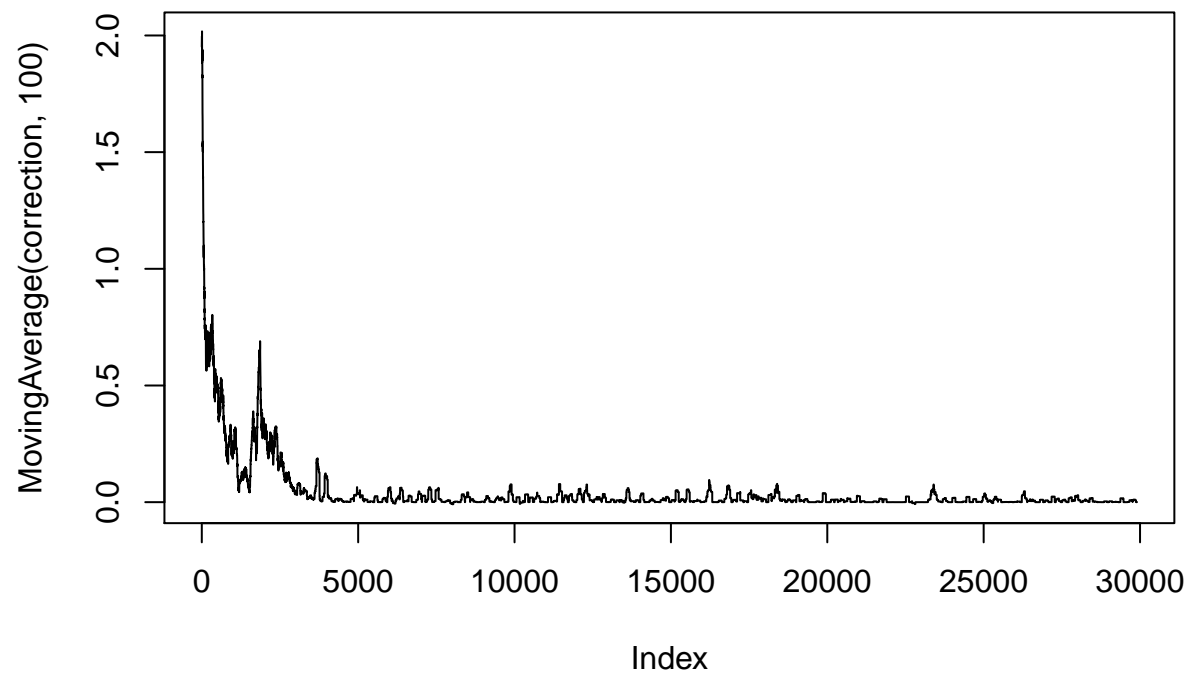  plot(MovingAverage(correction,100),type = "l")
```

```
}
```



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

# 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** | −0.19 / 0 ∨ 0.12 / 1.19 | −0.1 / 0.09 ∨ 0.16 / 1.58 | −0.34 / 0.32 > 2.1 / 0.57 | −0.27 / 0.64 ∨ 0.77 / 2.81 | −0.34 / 0 ∨ 0.19 / 3.75 | −0.1 / 0.41 < 0 / 0 | −0.1 / 0.03 < 0 / 0 | 0 / 0 ∨ 0 / 0 |
| **5** | 0.88 / 1.18 > 1.53 / 1.53 | 1.19 / 1.19 > 2.11 / 2.1 | 1.58 / 1.58 > 2.81 / 2.81 | 2.11 / 2.11 > 3.75 / 3.75 | 2.81 / 3.75 ∨ 2.83 / 5 | 0.11 / 1.75 < 0 / 0.2 | 0 / 0.28 < 0 / 0 | 0 / 0 ∧ 0 / 0 |
| **4** | 1.19 / 1.58 > 2.11 / 1.19 | 1.58 / 1.58 > 2.81 / 1.58 | 2.11 / 2.11 > 3.75 / 2.11 | 2.81 / 2.81 > 5 / 2.81 | 5 | 0 / 1.72 < 0 / 0 | 0 / 0 > 1 / 0 | 10 |
| **3** | 1.58 / 0.99 ∧ 1.3 / 0.69 | 1.37 / 2.02 > 2.11 / 0.36 | 2.81 / 1.52 ∧ 2.7 / 1.01 | 3.75 / 1.76 ∧ 2.42 / 1.6 | 4.45 / 0 ∧ 0 / 0 | 0 / 0 > 0 / 0 | 0 / 0 ∧ 0 / 0 | 0 / 0 ∧ 0 / 0 |
| **2** | 1.17 / 0.09 ∧ 0.01 / −0.27 | 0.16 / 0.71 < 0.1 / −0.19 | 0.4 / 0 > 1.8 / −0.34 | 2.79 / 0.1 ∧ 0 / −0.19 | 0 / 0.19 < 0 / −0.1 | 0 / 0 ∨ 0 / 0 | 0 / 0 < 0 / 0 | 0 / 0 > 0 / 0 |
| **1** | −1 | −1 | −1 | −1 | −1 | −1 | −1 | −1 |

x (vertical axis), y (horizontal axis)

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 )

The first thing we notice is that now the agent seems to find the +5 reward in each episode. This is because of a low value of $\epsilon = 0.1$ leading to less exploration and more trust in believing the existing policy. Therefore the episode correction converges faster as the solution is found faster. However, we notice that now generally there is a less chance of visiting +10 again due to lesser exploration. The gamma values don't really seem to influence the rewards much.

## 2.4 Environment C

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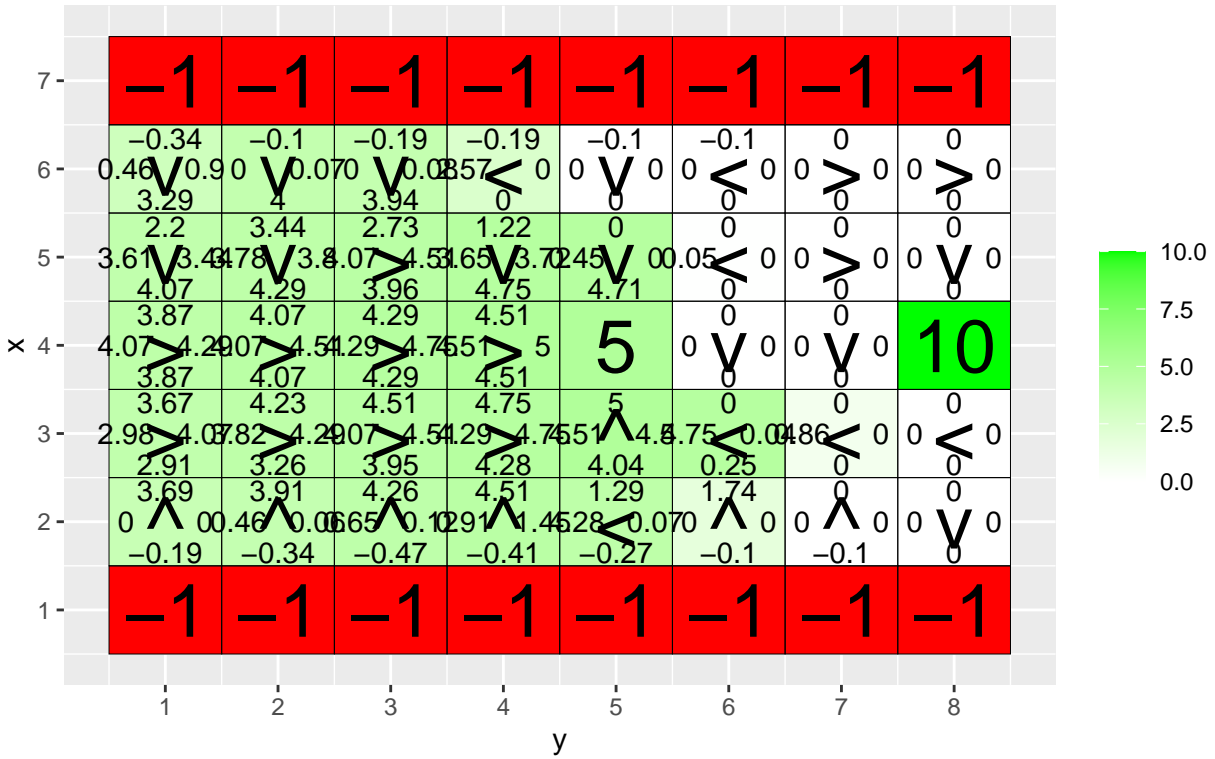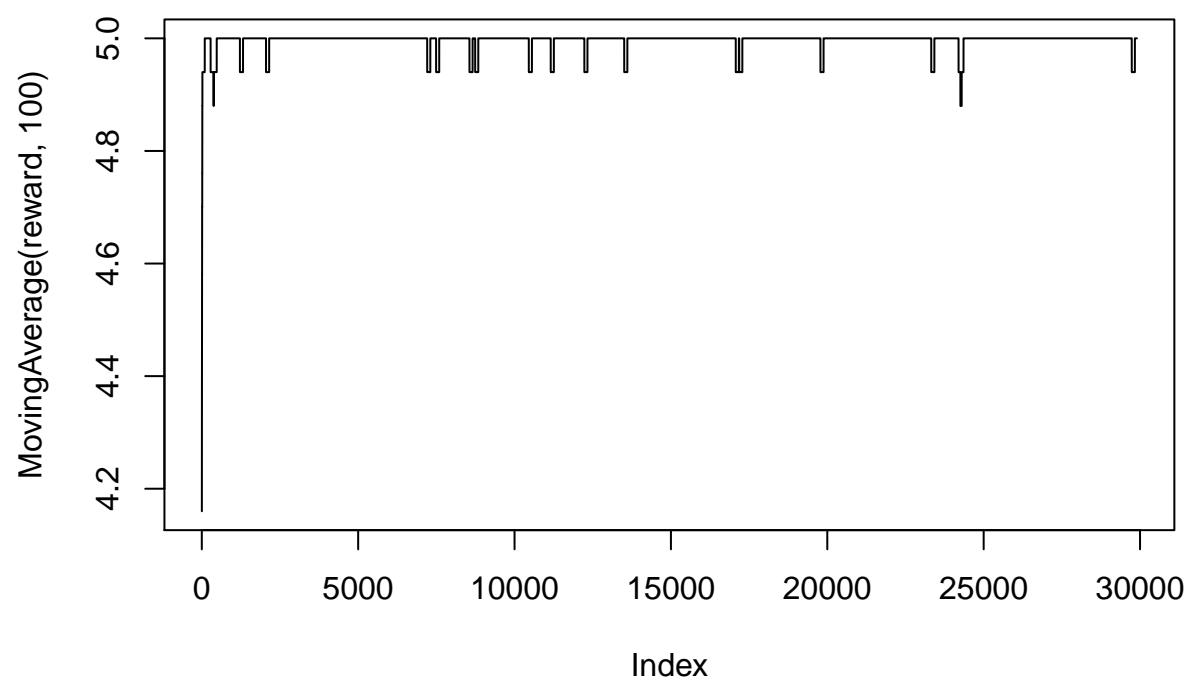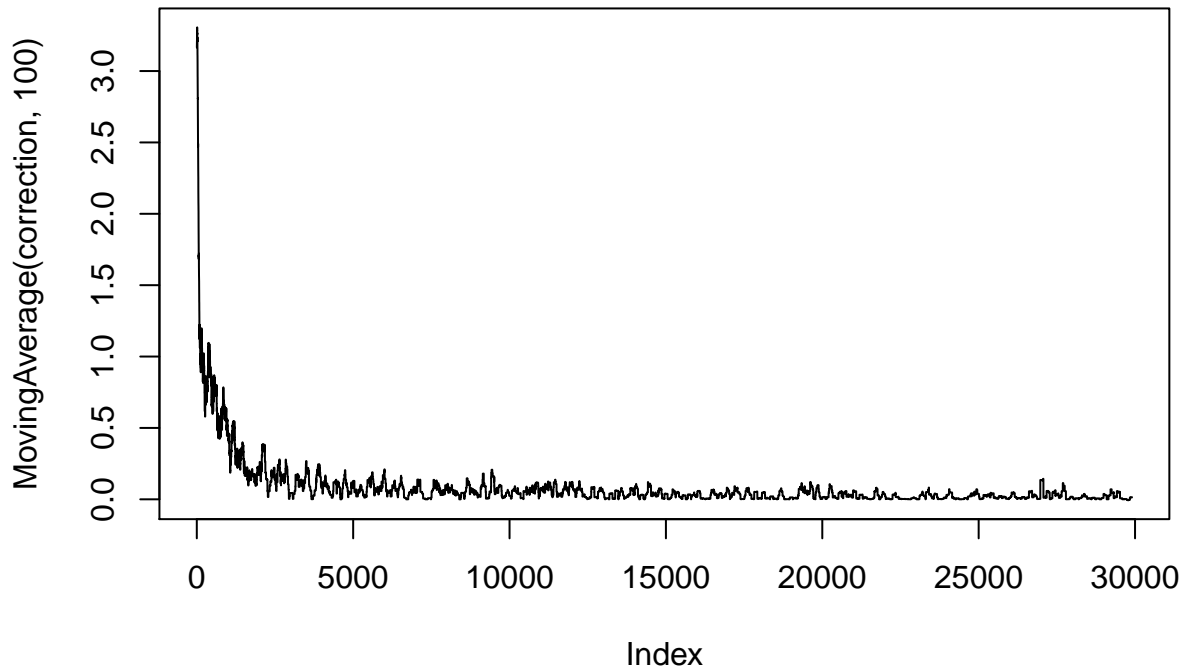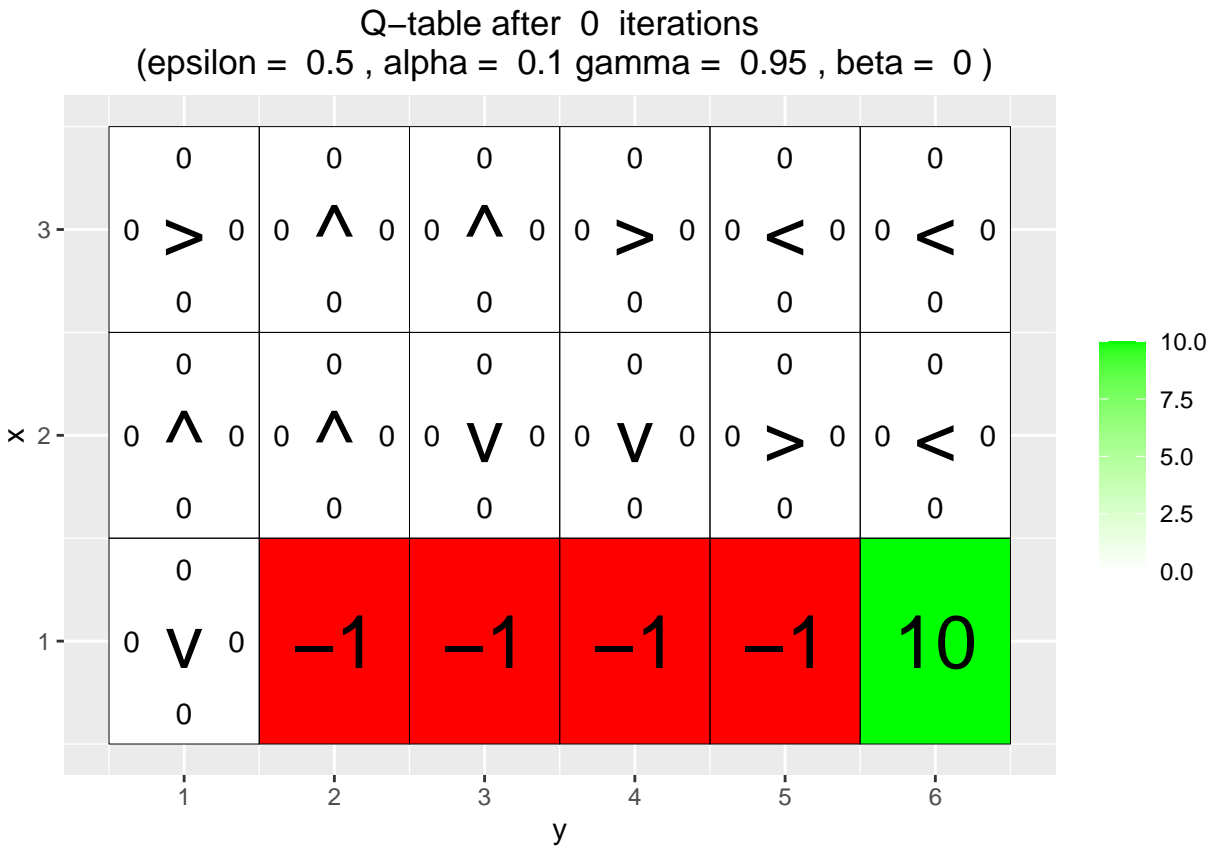# Environment C (the effect of beta).

H <- 3
W <- 6

reward_map <- matrix(0, nrow = H, ncol = W)
reward_map[1,2:5] <- -1
reward_map[1,6] <- 10
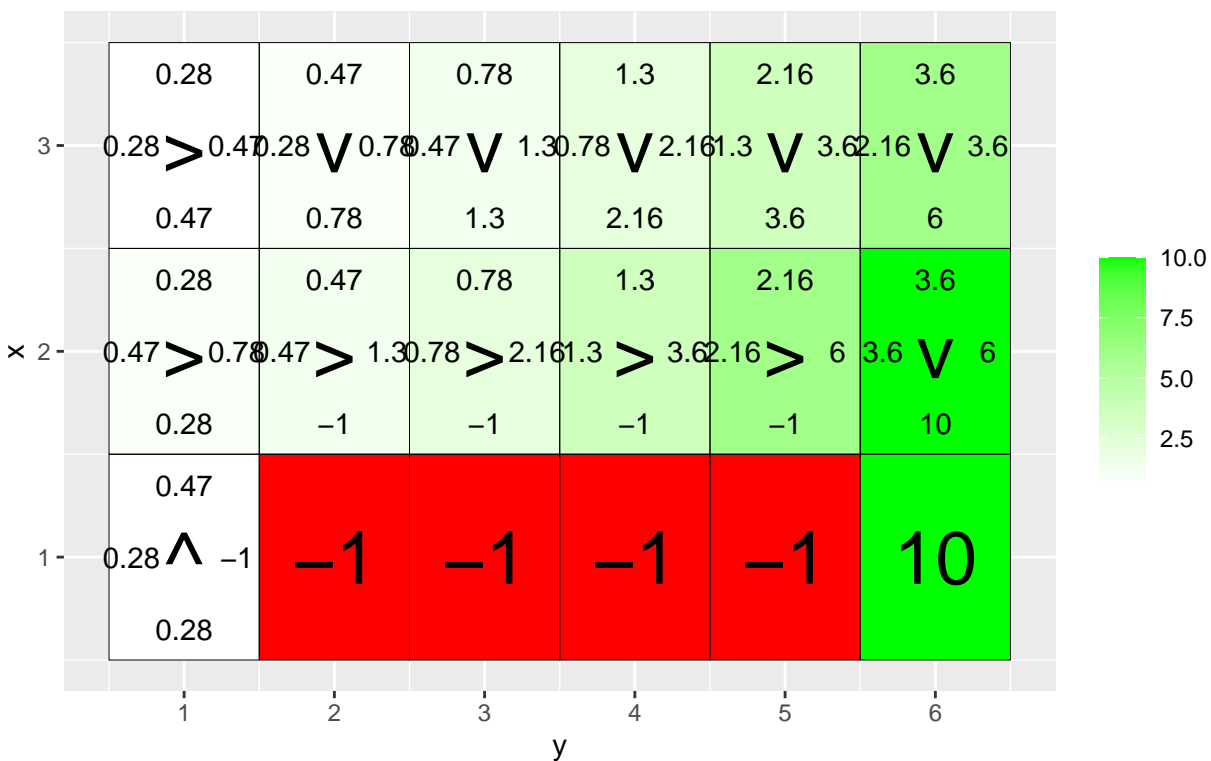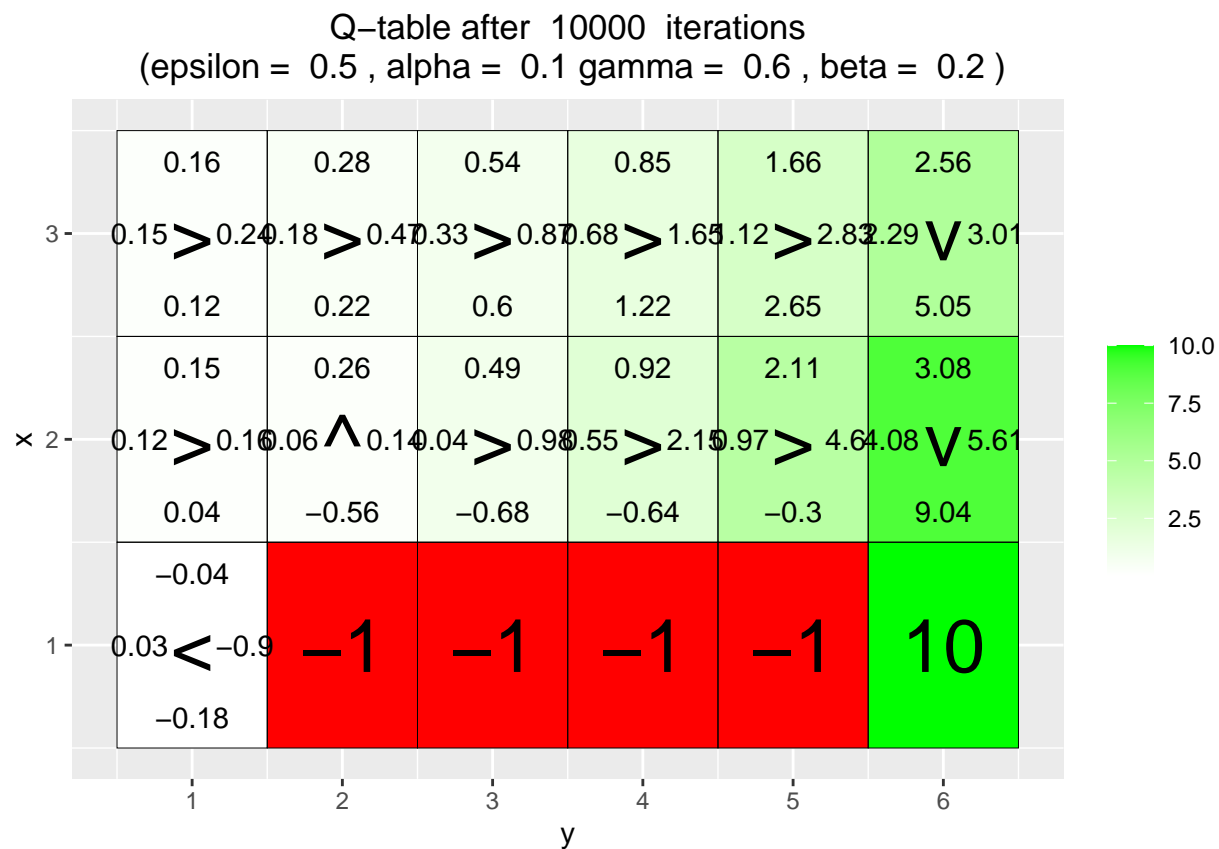
q_table <- array(0,dim = c(H,W,4))

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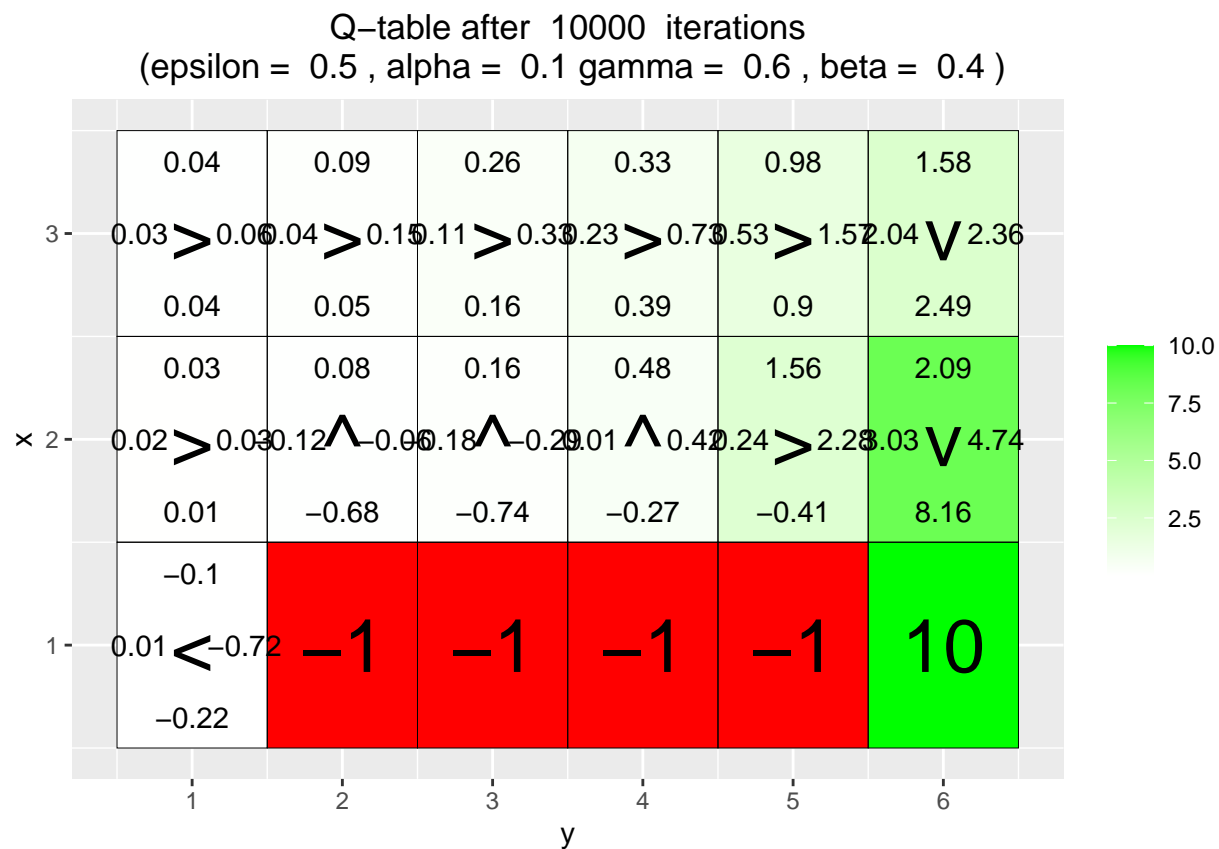
  for(i in 1:10000)
    foo <- q_learning(gamma = 0.6, beta = j, start_state = c(1,1))

  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 )

We see that as the slipping factor increases, the policy changes from finding the +10 reward to avoiding the -1 reward.

### 2.5 Environment D

1. Yes, we believe the agent has learned a good policy that is stochastic. Depending on where the goal position is, the policy is changed accordingly. It has done so because the policy parameters have updated to maximize the expected reward.

2. No Q-learning cannot be very easily applied here as there are a lot of state action pairs. There are n*(n-1) states (where n is the number of grids) and 4 actions. This would require updating a lot of Q values. Instead it would be easier to update the policy itself.

### 2.6 Environment E

1. No the agent has not learned a good policy because it was trained only on some goals at the side of the grid. These goals were not randomly distributed for training and made the agent learn a different policy (like going to the sides of the grid rather than trying to go closer to a goal)

2. Yes the results are different because both these environments had different kinds of training -where one was biased and one was not.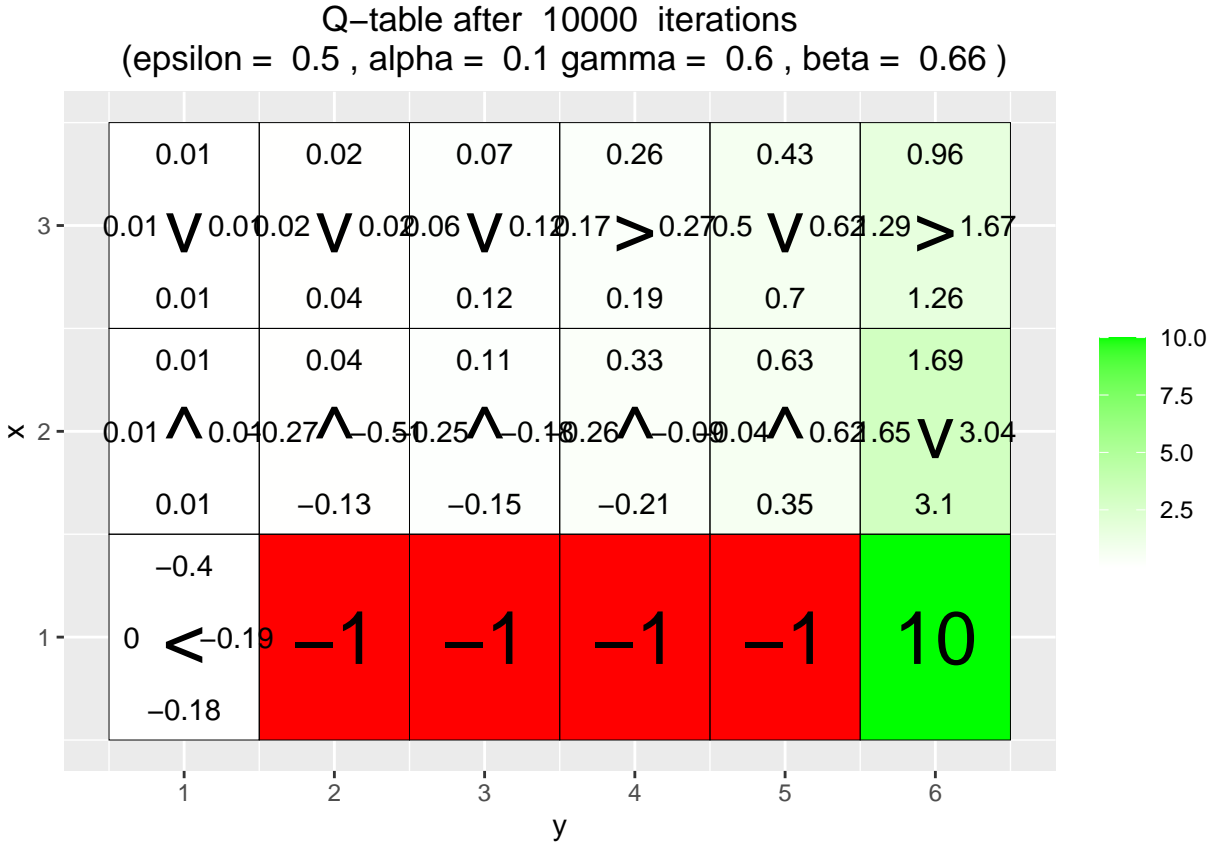