

LAB3 Hjalmar Oehman

Hjalmar Öhman

2024-10-03

3.1

- You are provided with a gridworld environment where an agent can move up, down, left, or right. The agent starts each episode at a specific location and receives rewards based on the environment's layout.
- Implement the **Greedy** and **Epsilon-Greedy** policies in the functions `GreedyPolicy` and `EpsilonGreedyPolicy` respectively.
- Implement the **Q-learning** algorithm in the function `q_learning`, which will update the Q-table based on the agent's actions and rewards.
- Run 10000 episodes of Q-learning with $\epsilon = 0.5$, $\beta = 0$, $\alpha = 0.1$, and $\gamma = 0.95$. Visualize the results at various stages (after 10, 100, 1000, and 10000 episodes).
- Answer the following questions:
 1. What has the agent learned after the first 10 episodes?
 2. Is the final greedy policy (after 10000 episodes) optimal for all states? Why or why not?
 3. Do the learned values in the Q-table reflect the existence of multiple paths to the positive reward?

```
rm(list = ls())
# 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)
set.seed(1234)

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))
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.
  q_values = q_table[x, y, ]

```

```

best_actions = which(q_values == max(q_values))
if (length(best_actions) > 1) {
  action = sample(best_actions, 1)
} else {
  action = best_actions
}
return (action)
}

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}.

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

  return (action)
}

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){

#Initialize
current_state = start_state
episode_correction = 0
reward = 0

repeat{
    # Current state
    x = current_state[1]
    y = current_state[2]

    # Action
    action = EpsilonGreedyPolicy(x, y, epsilon)

    # Update state
    next_state = transition_model(x, y, action, beta)
    next_x = next_state[1]
    next_y = next_state[2]

    # Reward
    reward = reward_map[next_x, next_y]

    # New max Q value
    max_q_next = max(q_table[next_x, next_y, ])

    # Correction
    correction = reward + gamma * max_q_next - q_table[x, y, action]

    # Update q_table based on correction
    q_table[x, y, action] <- q_table[x, y, action] + alpha * correction

    # Accumulate corrections
    episode_correction = episode_correction + correction

    # Update state
    current_state = next_state

    # End the episode if a terminal state (non-zero reward) is reached
    if (reward != 0) {
        break
    }
}

return (c(reward, episode_correction))
}

```

3.2

Environment A: - Set up a 5x7 grid with rewards at specific states and analyze the agent's learning process by visualizing its Q-table and policies at different stages of learning.

```

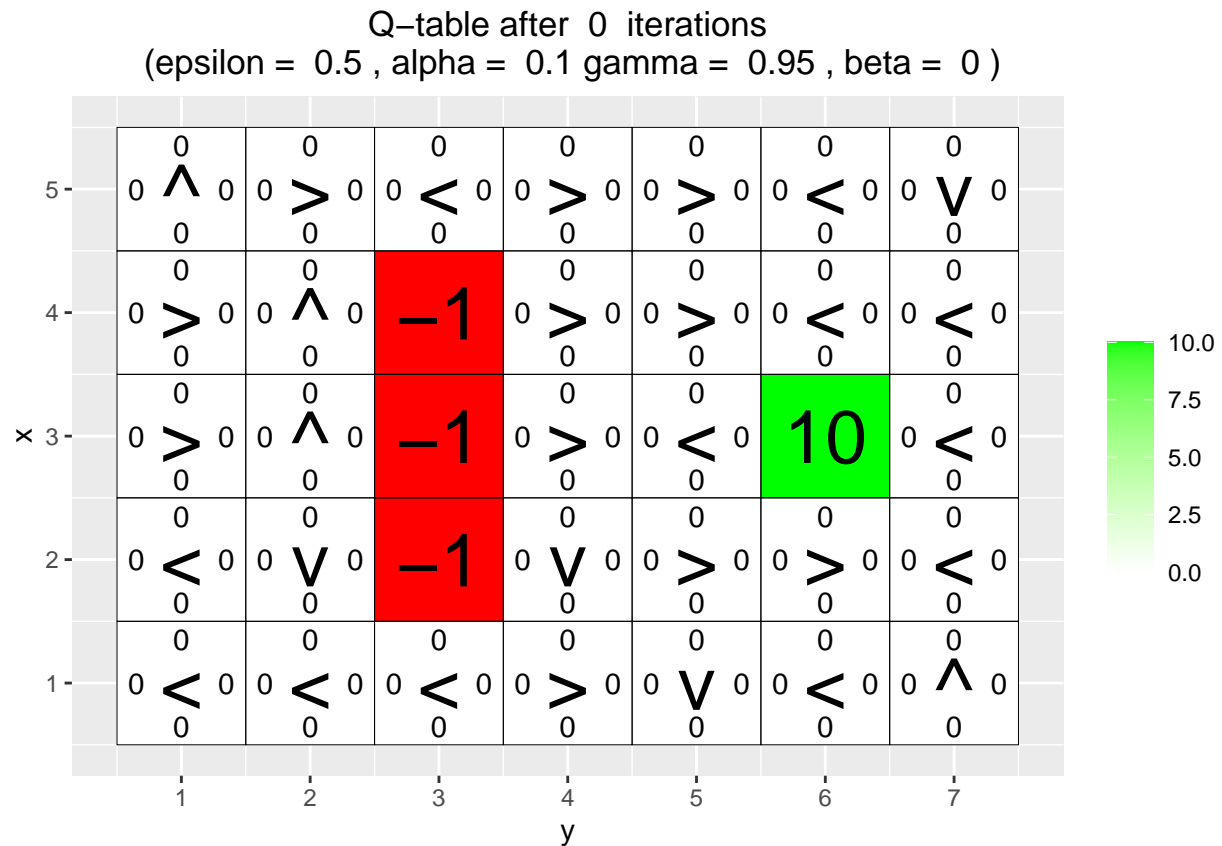
# Environment A (learning)
set.seed(1234)
H <- 5
W <- 7

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()

```



```

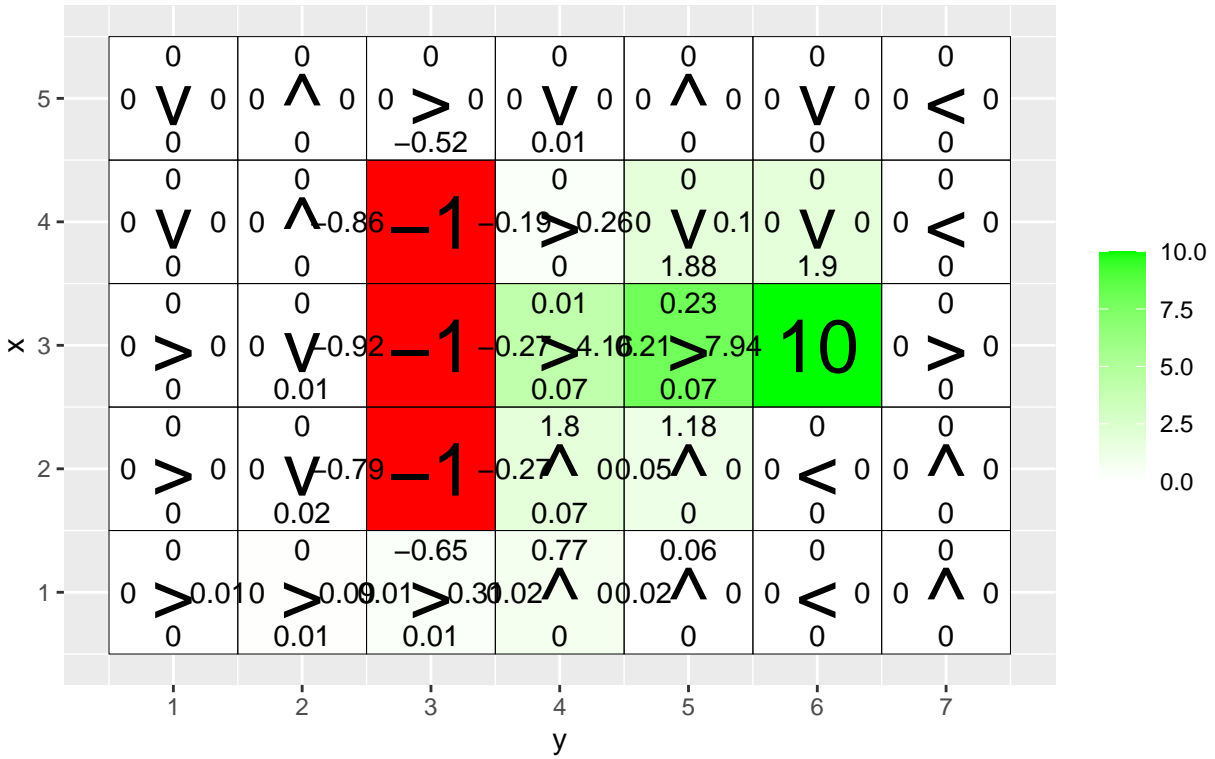
for(i in 1:10000){
  foo <- q_learning(start_state = c(3,1))

  if(any(i==c(10,100,1000,10000)))
    vis_environment(i)
}

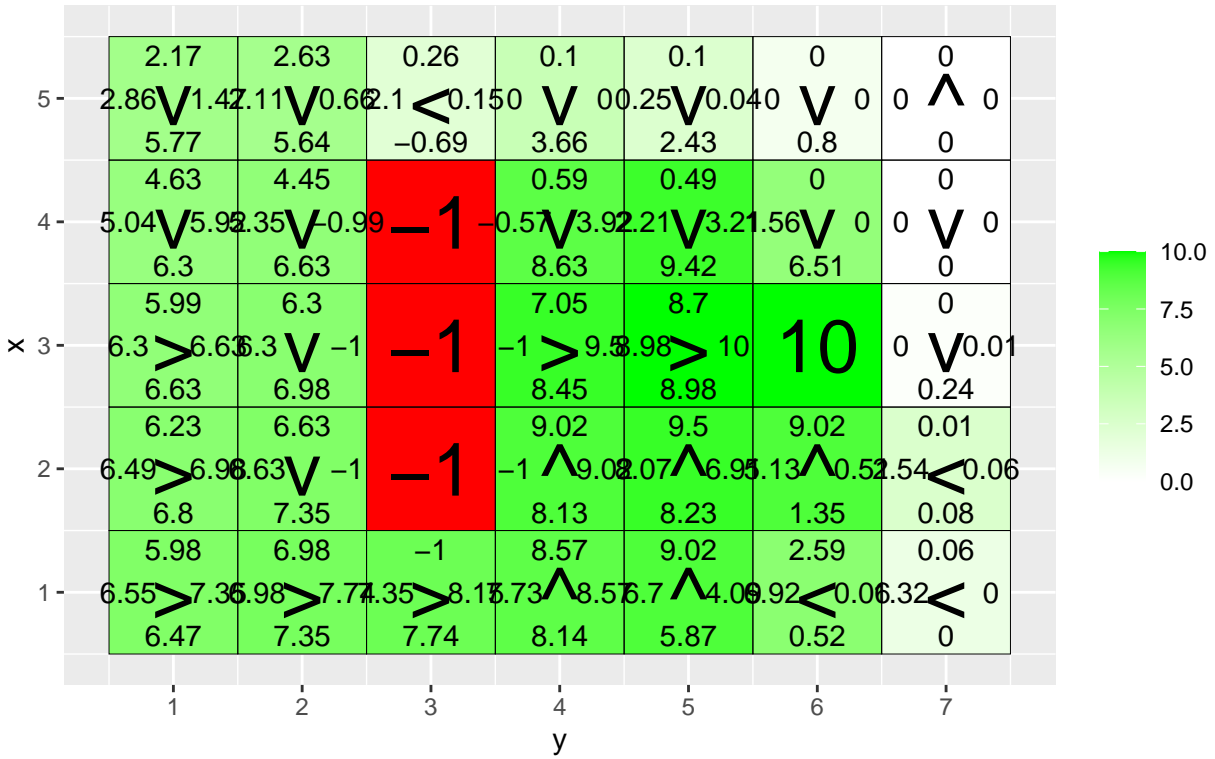
```

[illegible]

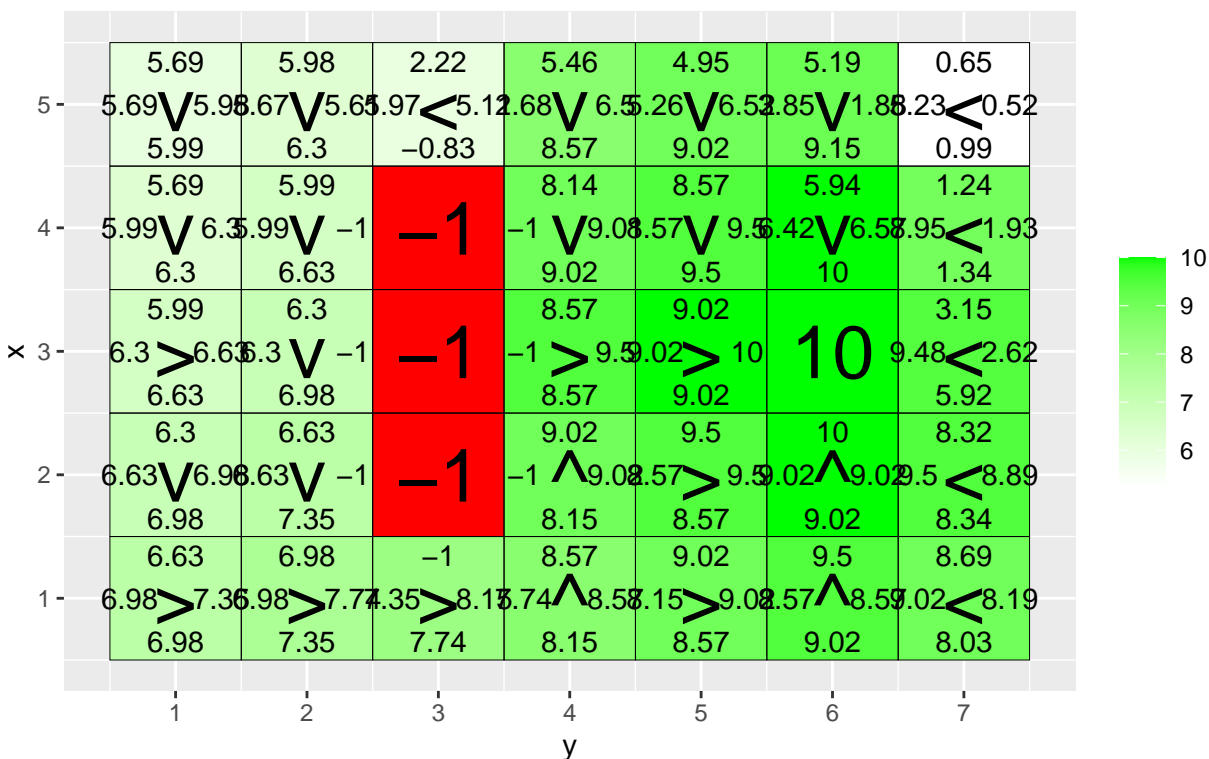
Q-table after 100 iterations
 (epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0)



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)



Q: What has the agent learned after the first 10 episodes ? A: To not enter(4, 3)

Q: Is the final greedy policy (after 10000 episodes) optimal for all states, i.e. not only for the initial state ? Why / Why not ? A: No. For example, the action policy for state (5, 3) is to go left, but the optimal action is right. This is due to the agent not having had the chance of exploring this state enough times.

Q: 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 ? A: No, it always tries to go below. We could either enforce a greater epsilon, or change start_state to be further up (e.g (1, 5))

3.3

Environment B: - Investigate the effects of different values of ϵ and γ by running 30000 episodes of Q-learning with various configurations.

```
# Environment B (the effect of epsilon and gamma)
set.seed(1234)

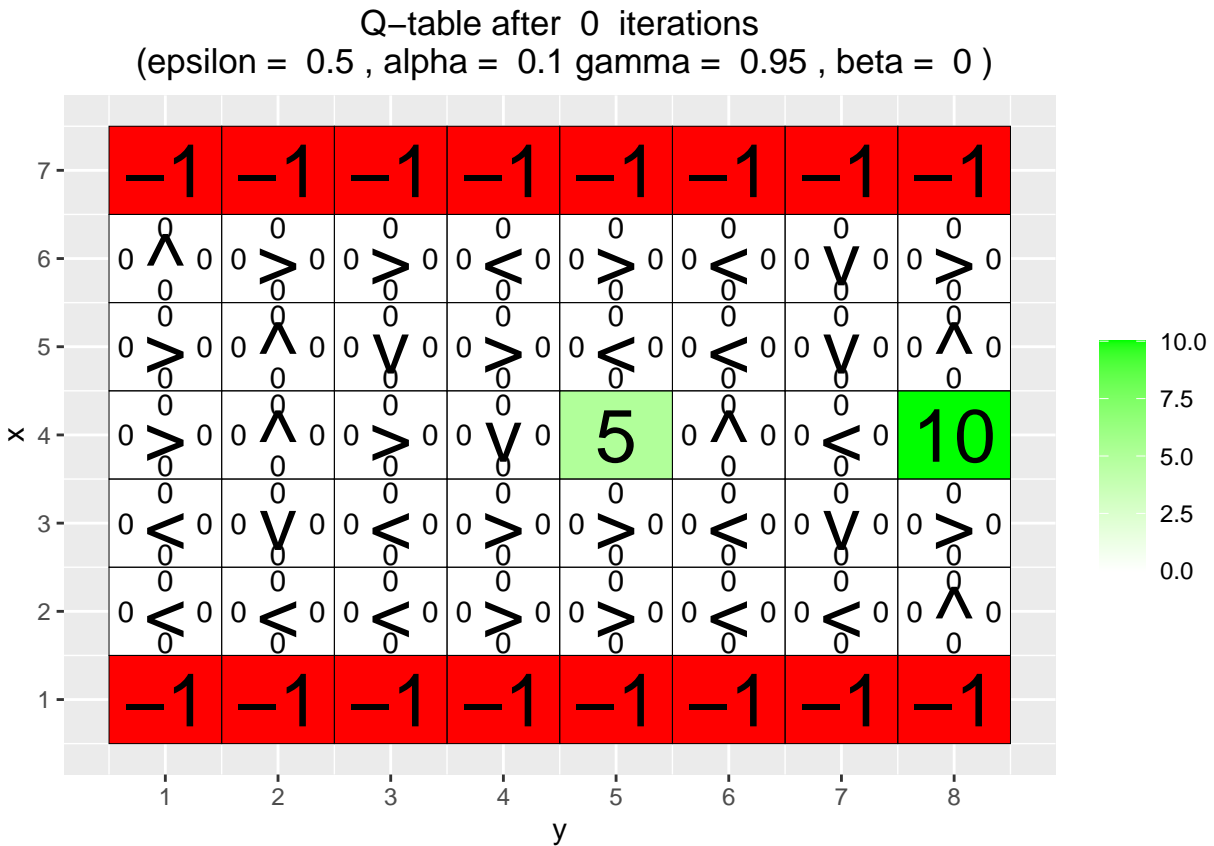
H <- 7
W <- 8

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()
```



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

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

  return (rsum)
}

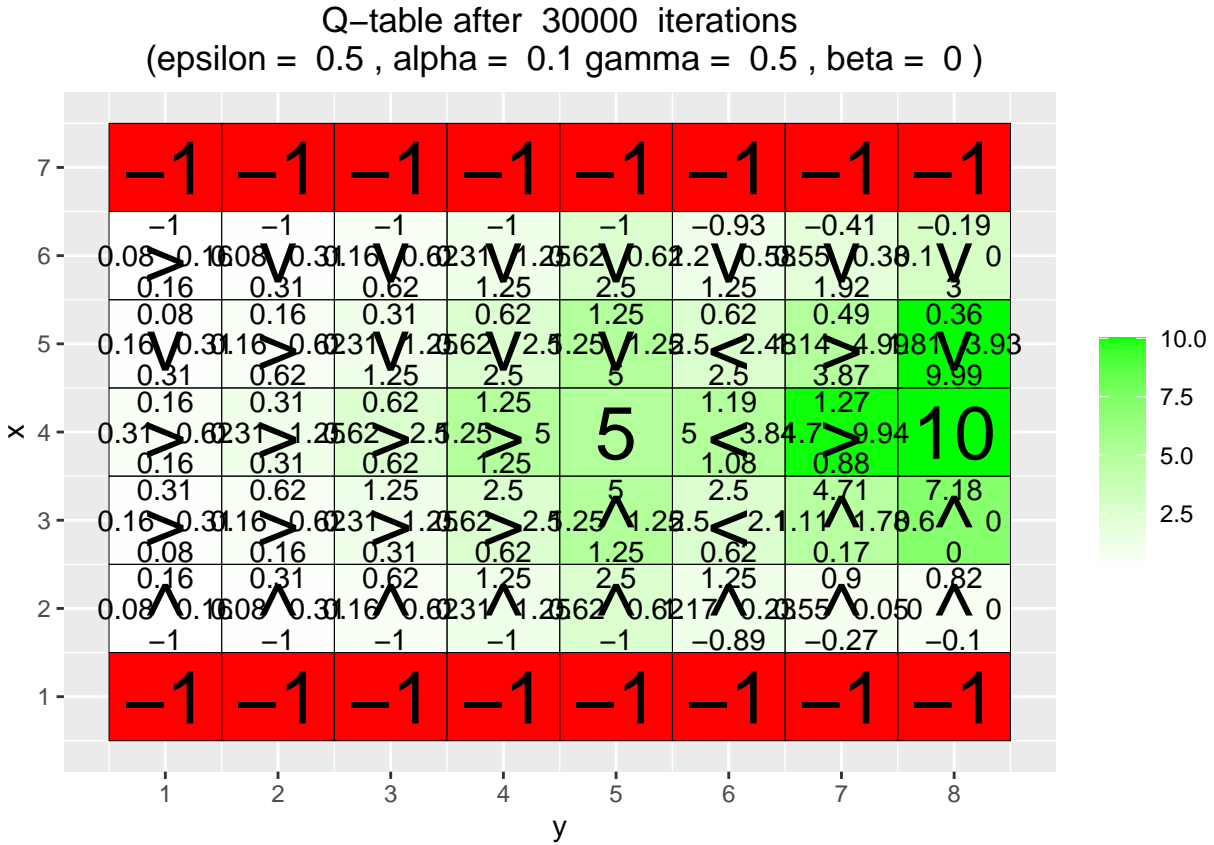
for(j in c(0.5,0.75,0.95)){
  q_table <- array(0,dim = c(H,W,4))
  reward <- NULL
  correction <- NULL

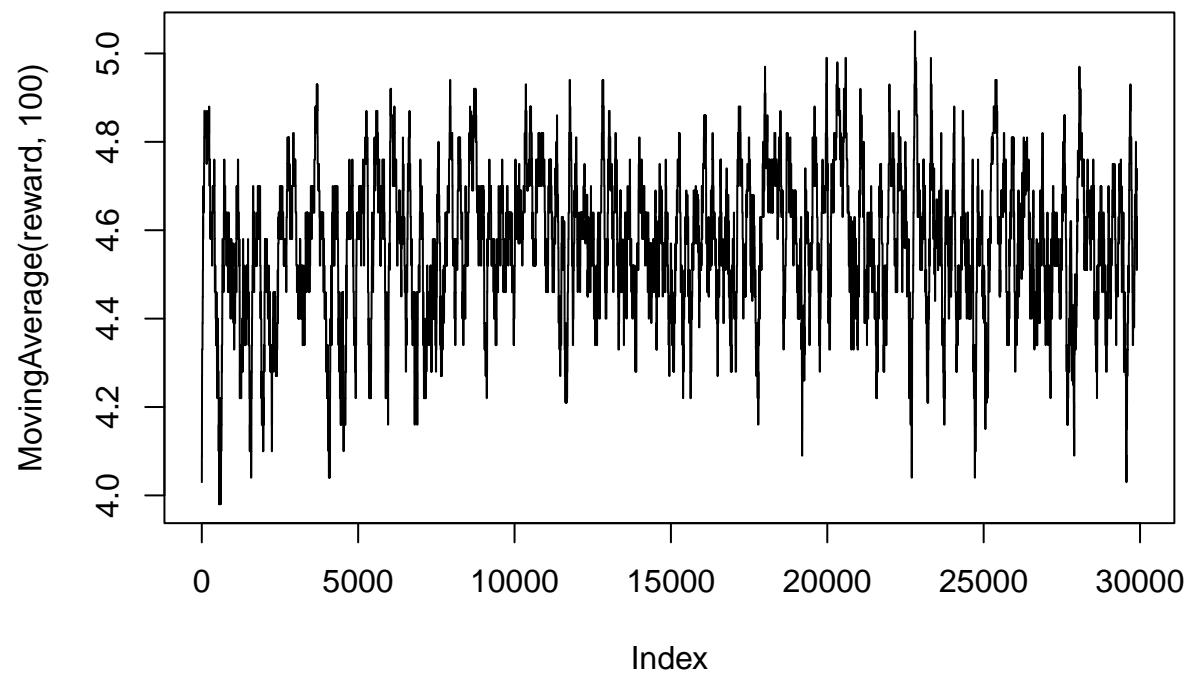
  for(i in 1:30000){
    foo <- q_learning(gamma = j, start_state = c(4,1))
    reward <- c(reward,foo[1])
    correction <- c(correction,foo[2])
  }
}
```

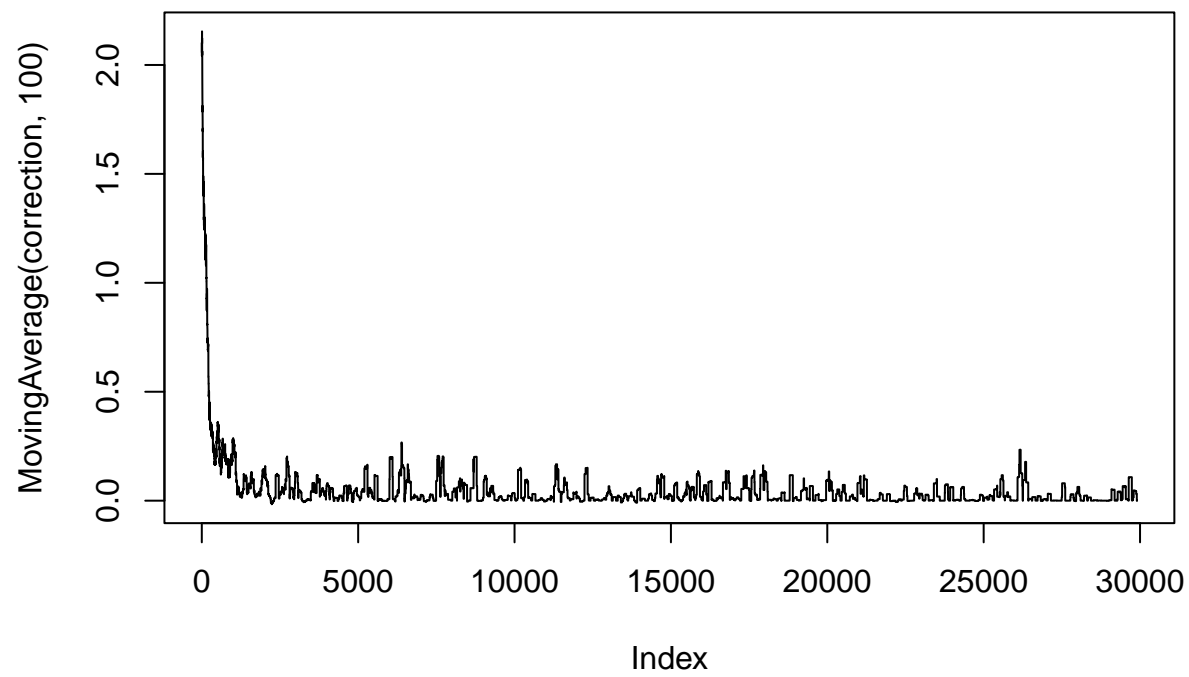
```

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

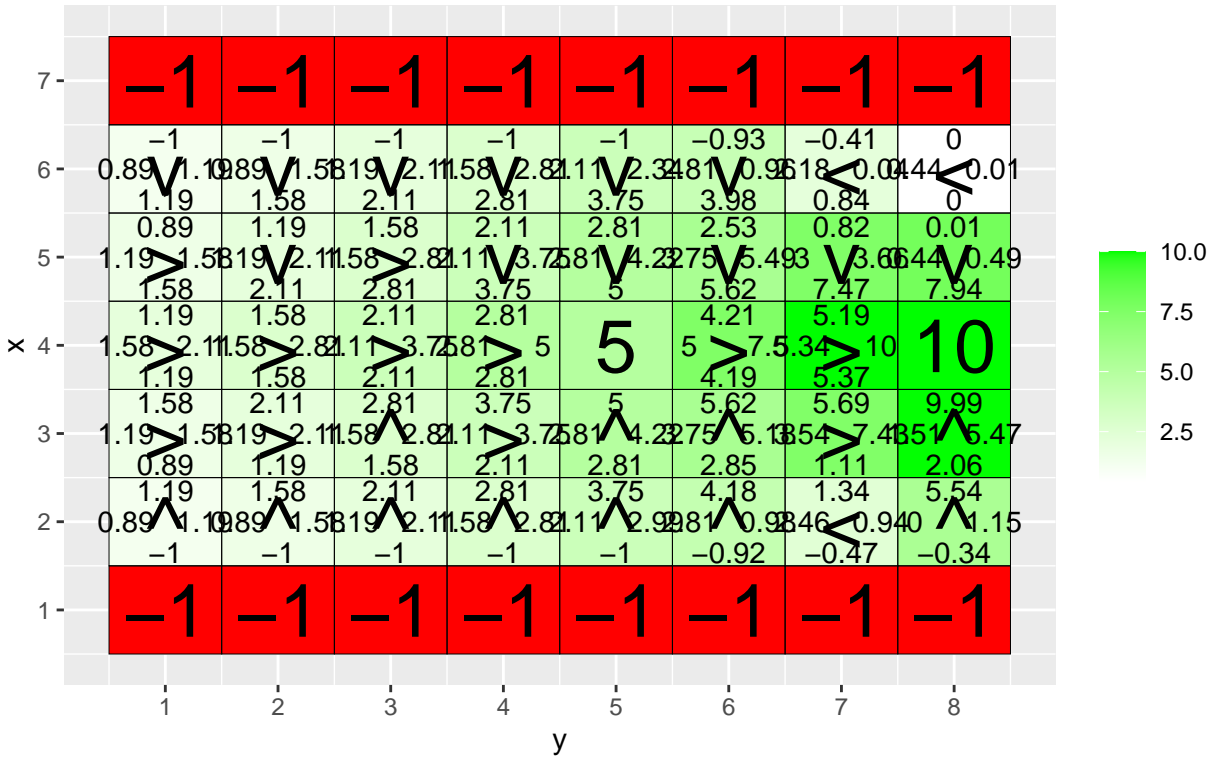
```

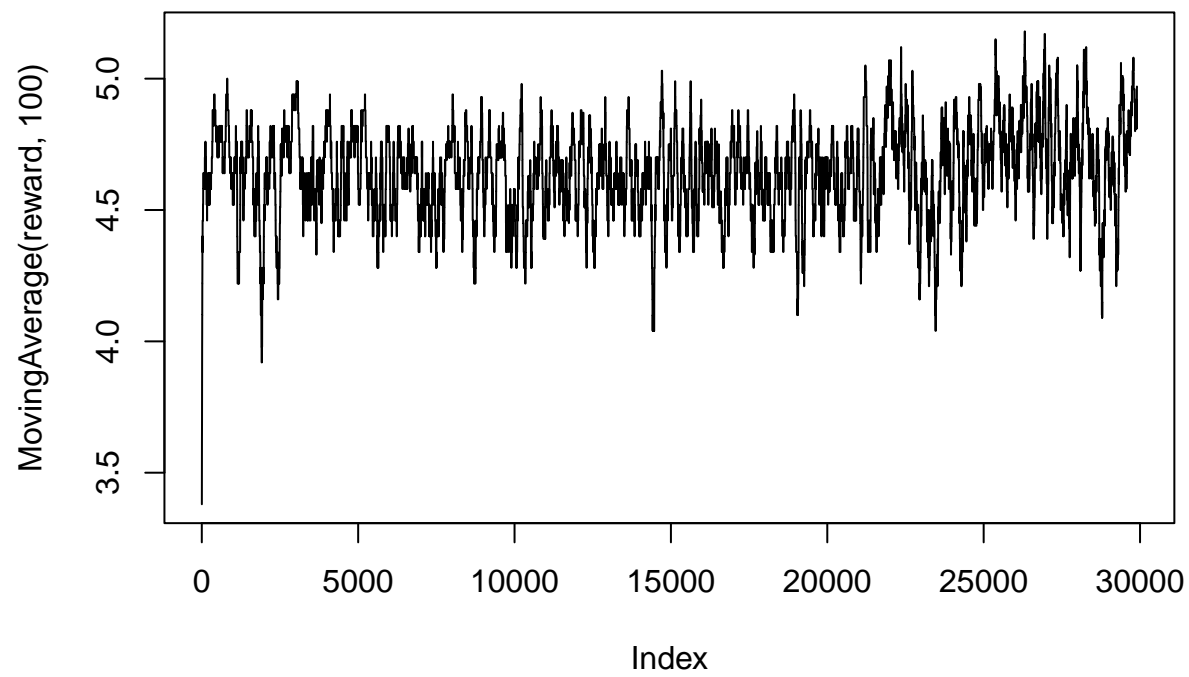


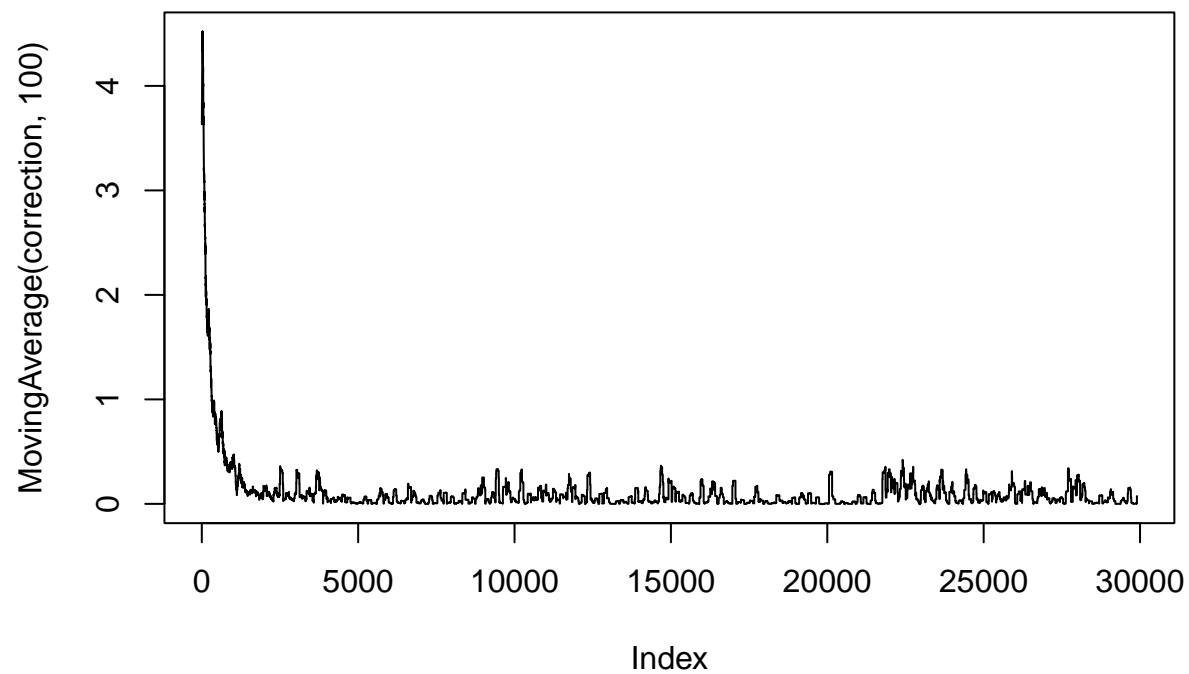




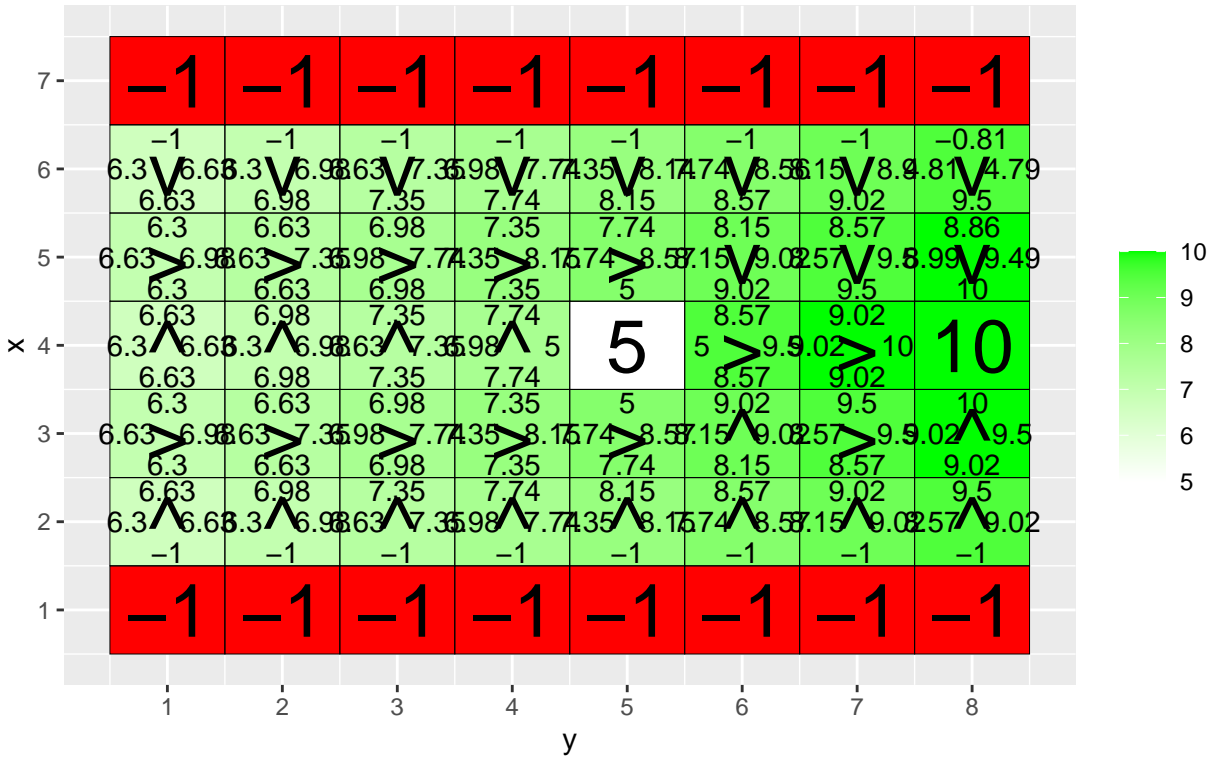
Q-table after 30000 iterations
(epsilon = 0.5 , alpha = 0.1 gamma = 0.75 , beta = 0)

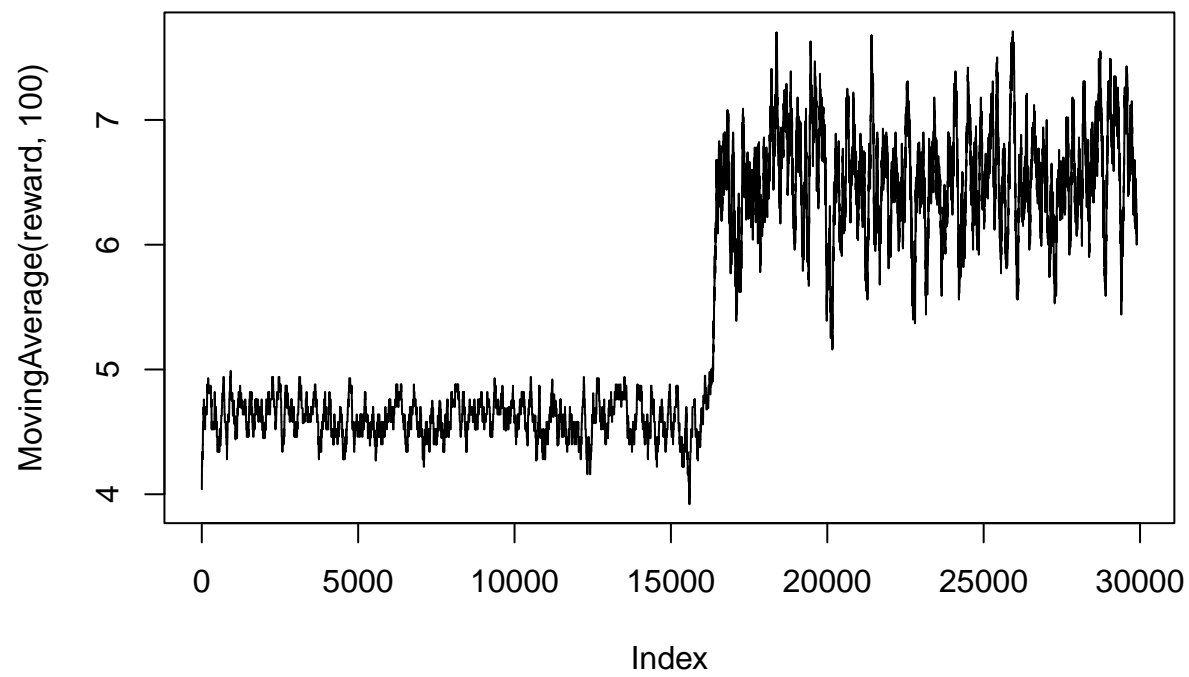


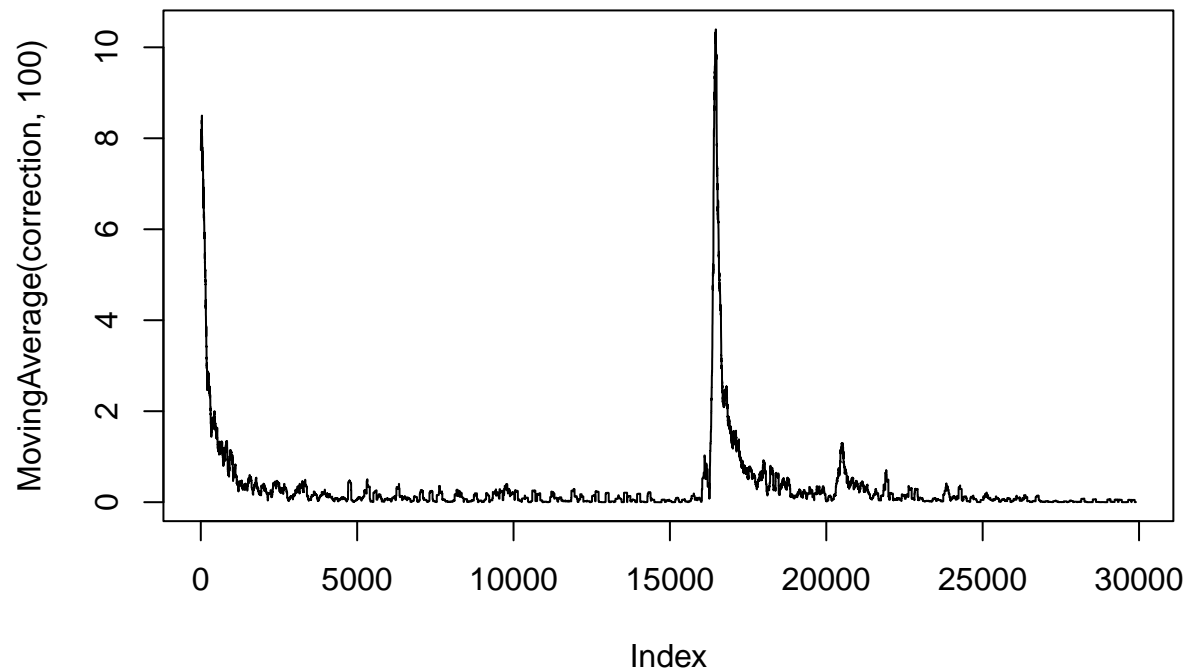




Q-table after 30000 iterations
 (epsilon = 0.5 , alpha = 0.1 gamma = 0.95 , beta = 0)





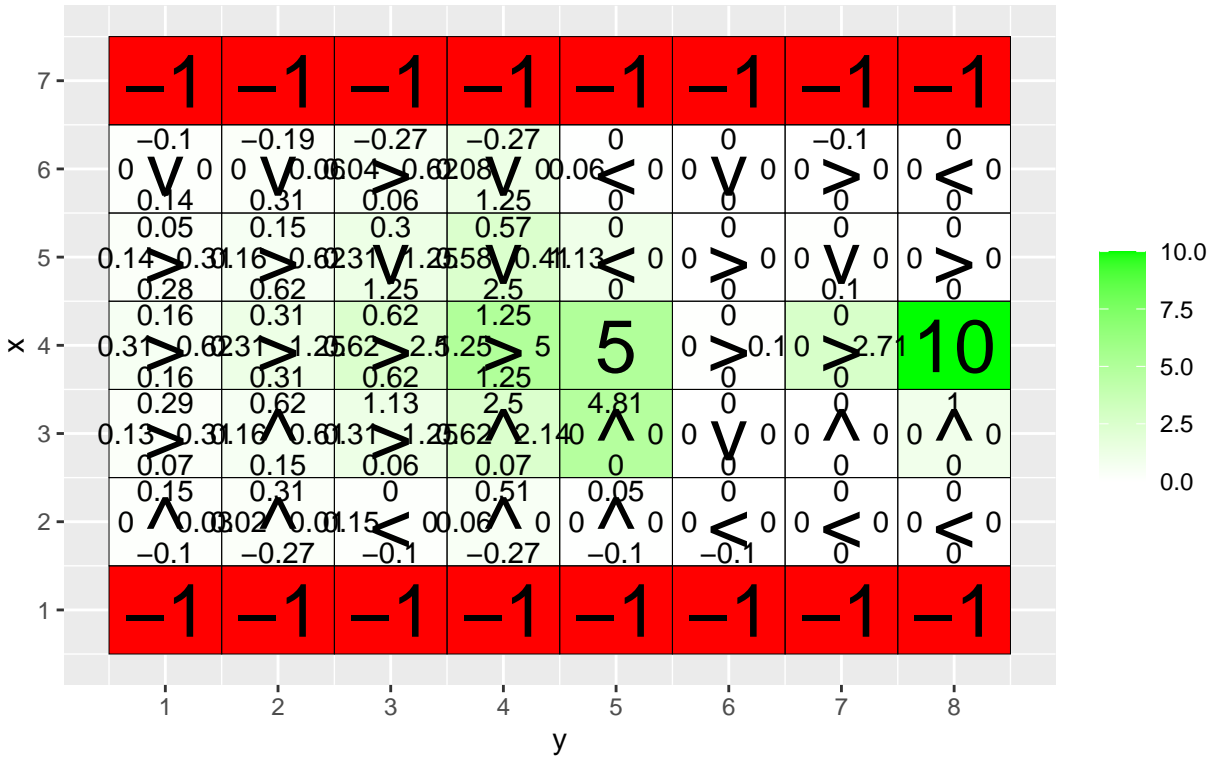


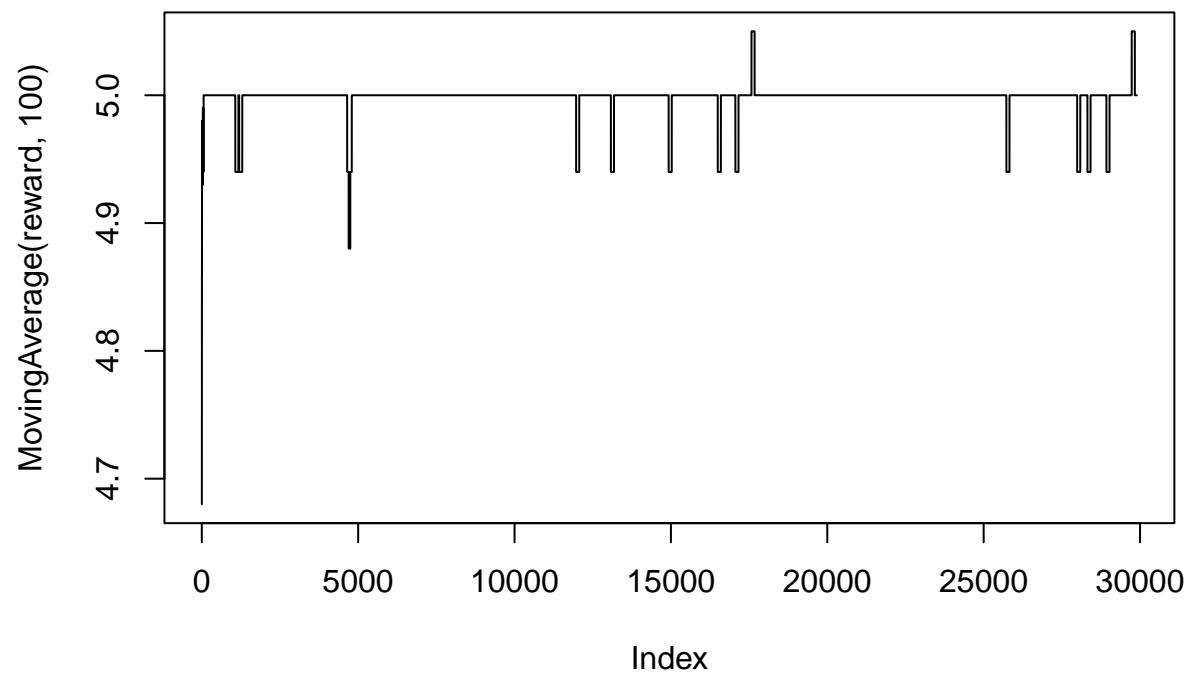
```
for(j in c(0.5,0.75,0.95)){
  q_table <- array(0,dim = c(H,W,4))
  reward <- NULL
  correction <- NULL

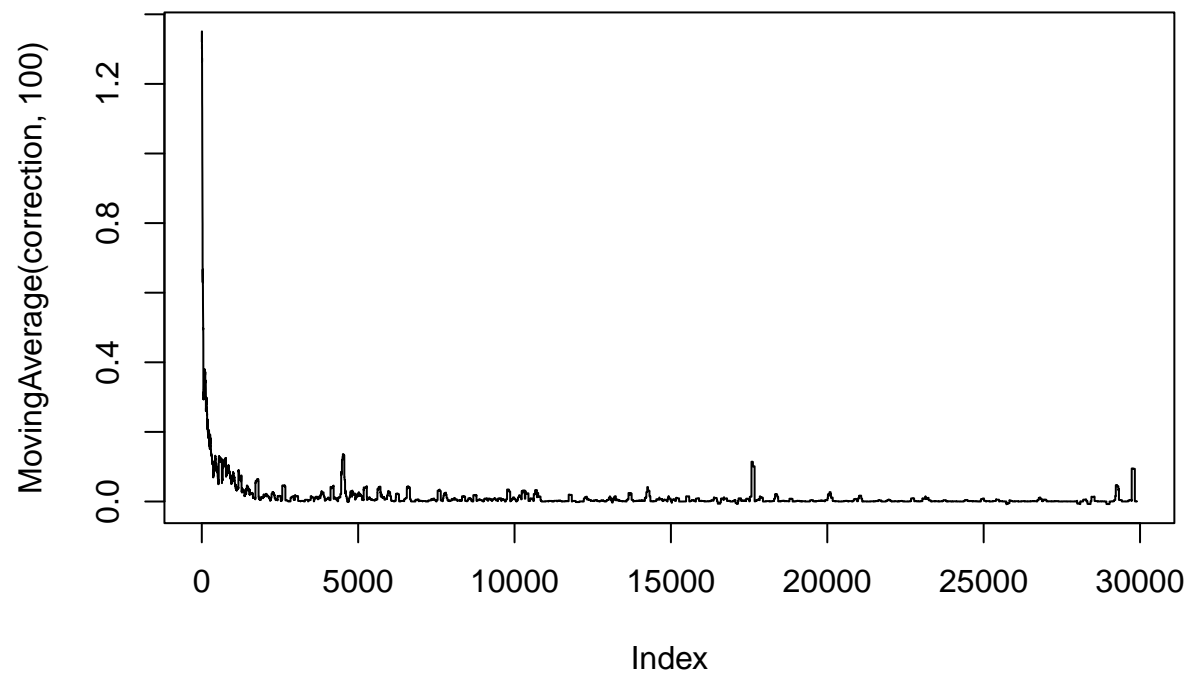
  for(i in 1:30000){
    foo <- q_learning(epsilon = 0.1, gamma = j, start_state = c(4,1))
    reward <- c(reward,foo[1])
    correction <- c(correction,foo[2])
  }

  vis_environment(i, epsilon = 0.1, gamma = j)
  plot(MovingAverage(reward,100),type = "l")
  plot(MovingAverage(correction,100),type = "l")
}
```

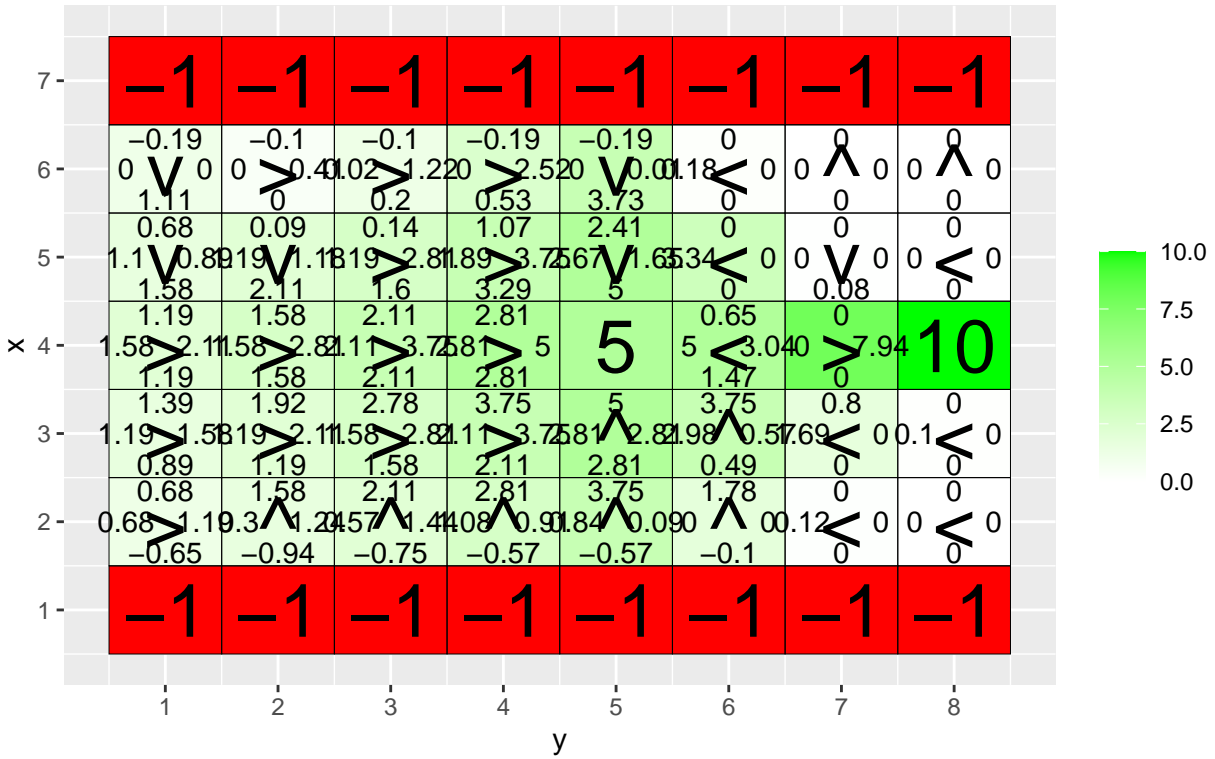
Q-table after 30000 iterations
(epsilon = 0.1 , alpha = 0.1 gamma = 0.5 , beta = 0)

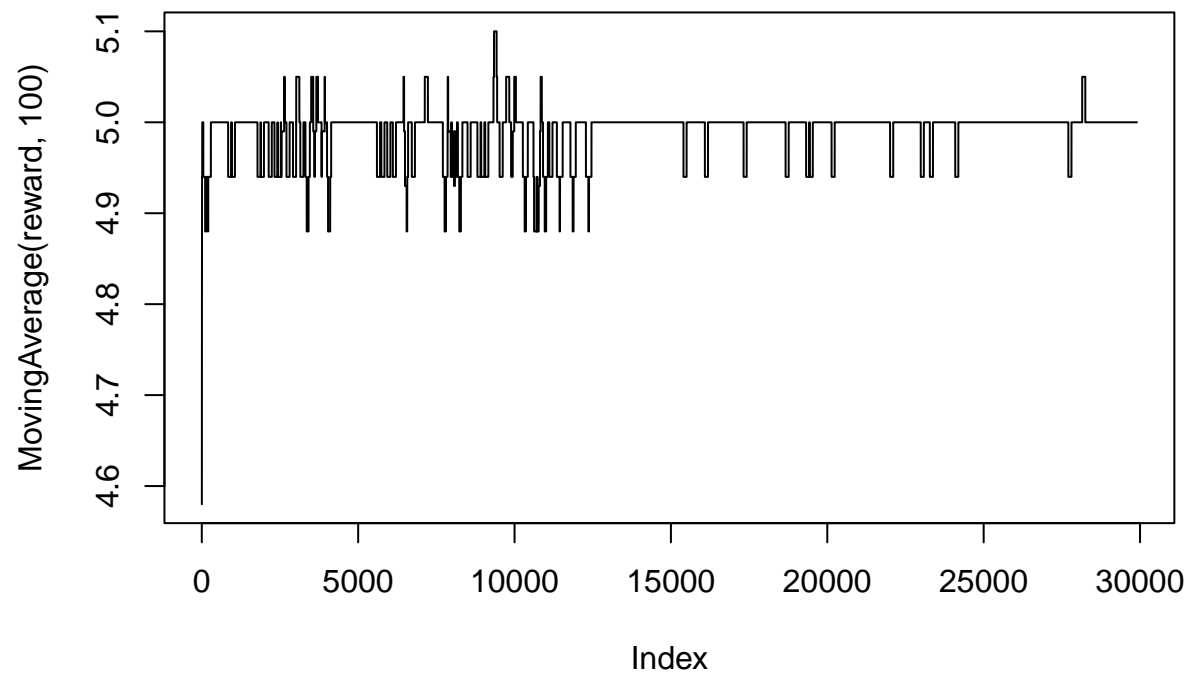


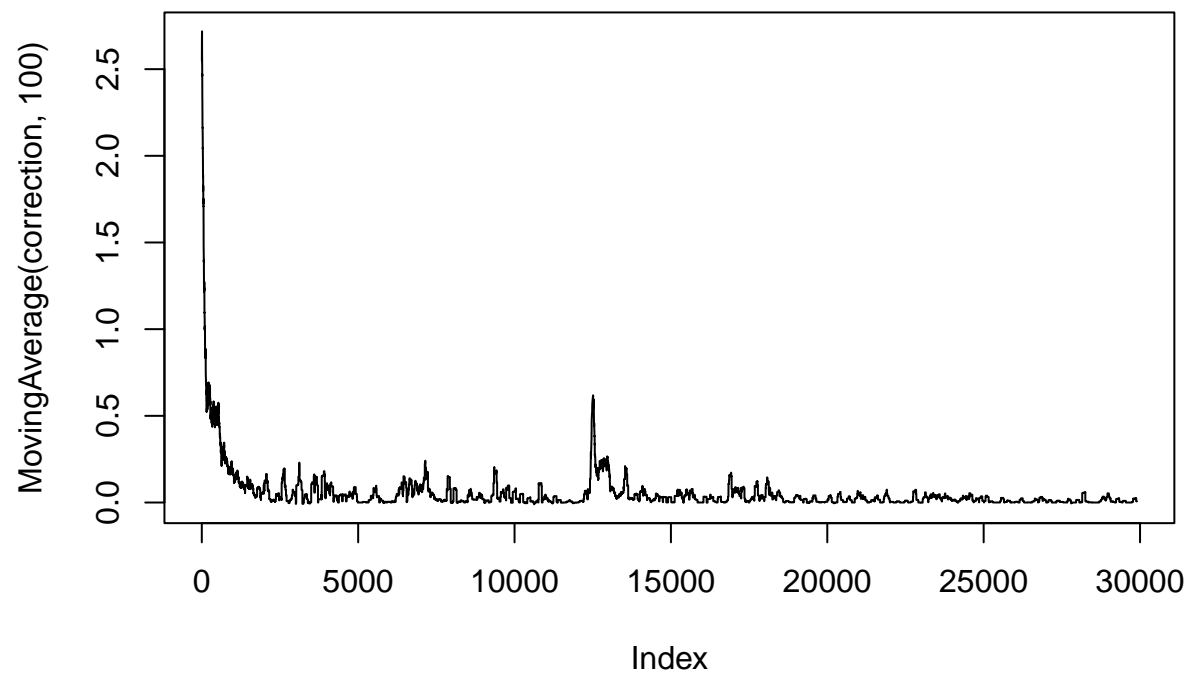




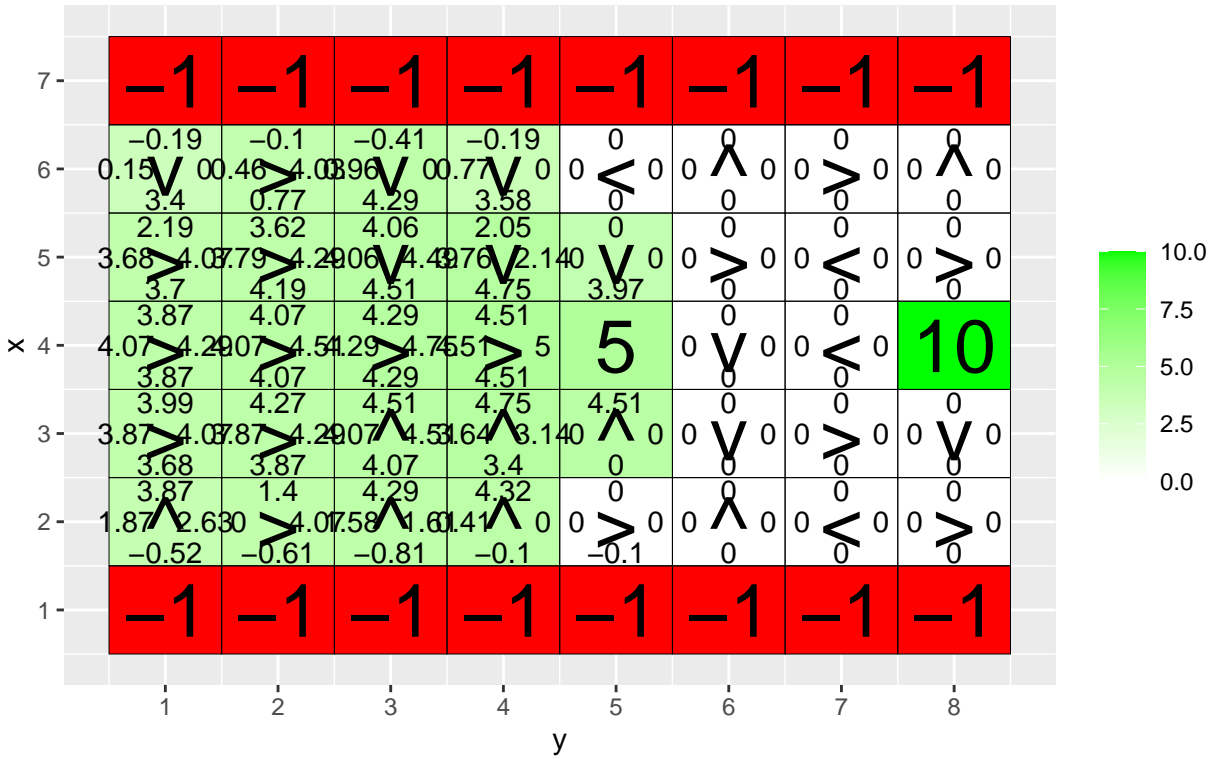
Q-table after 30000 iterations
(epsilon = 0.1 , alpha = 0.1 gamma = 0.75 , beta = 0)

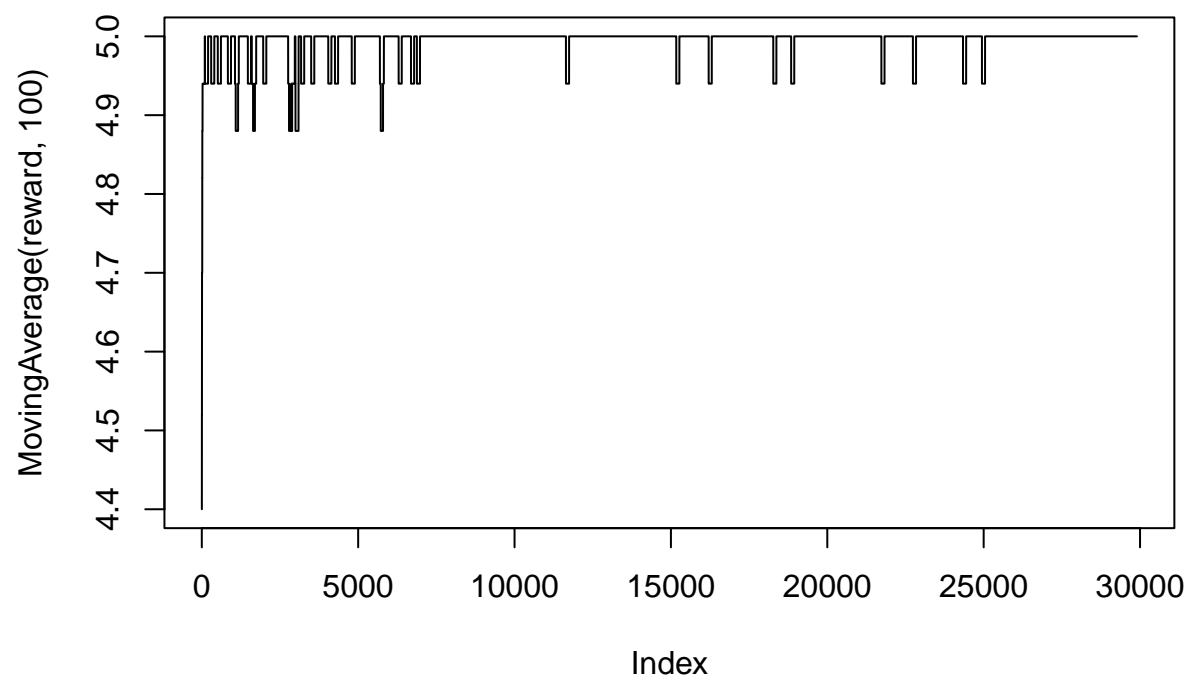


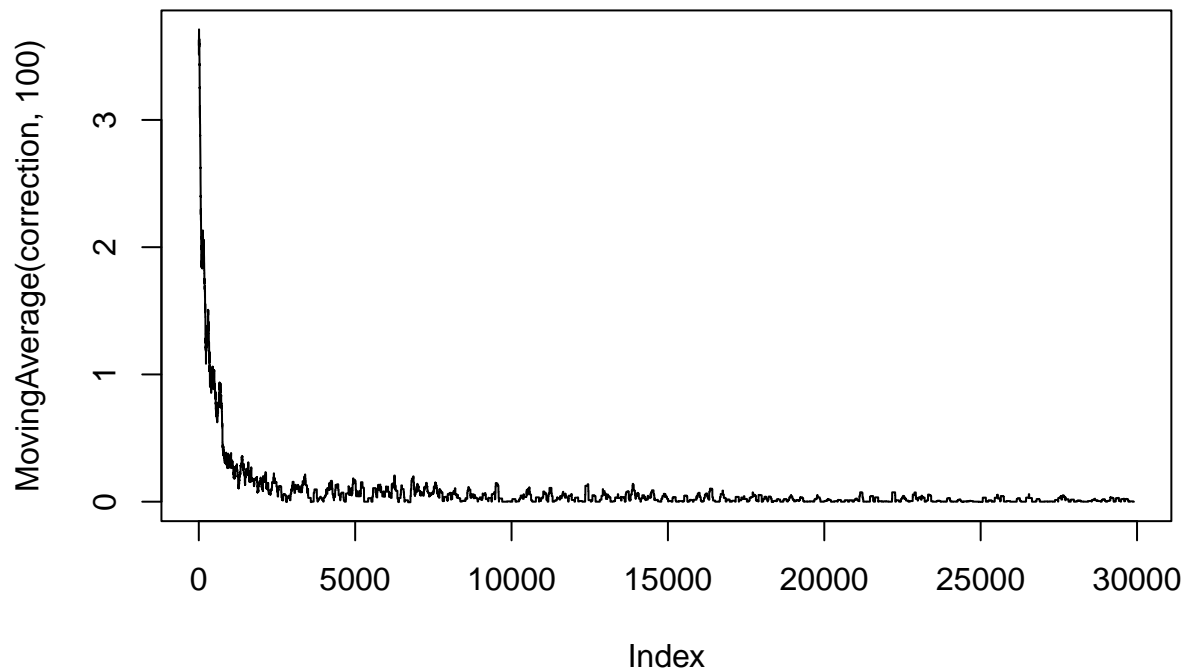




Q-table after 30000 iterations
(epsilon = 0.1 , alpha = 0.1 gamma = 0.95 , beta = 0)







Q: Investigate how the ϵ and γ parameters affect the learned policy A: A high γ changes the action policy to be more prone going towards 10 over 5. A low ϵ leads to the agent not exploring certain states.

3.4

Environment C: - Investigate how the β parameter (which controls the probability of slipping) affects the learned policy by running 10000 episodes with different values of β .

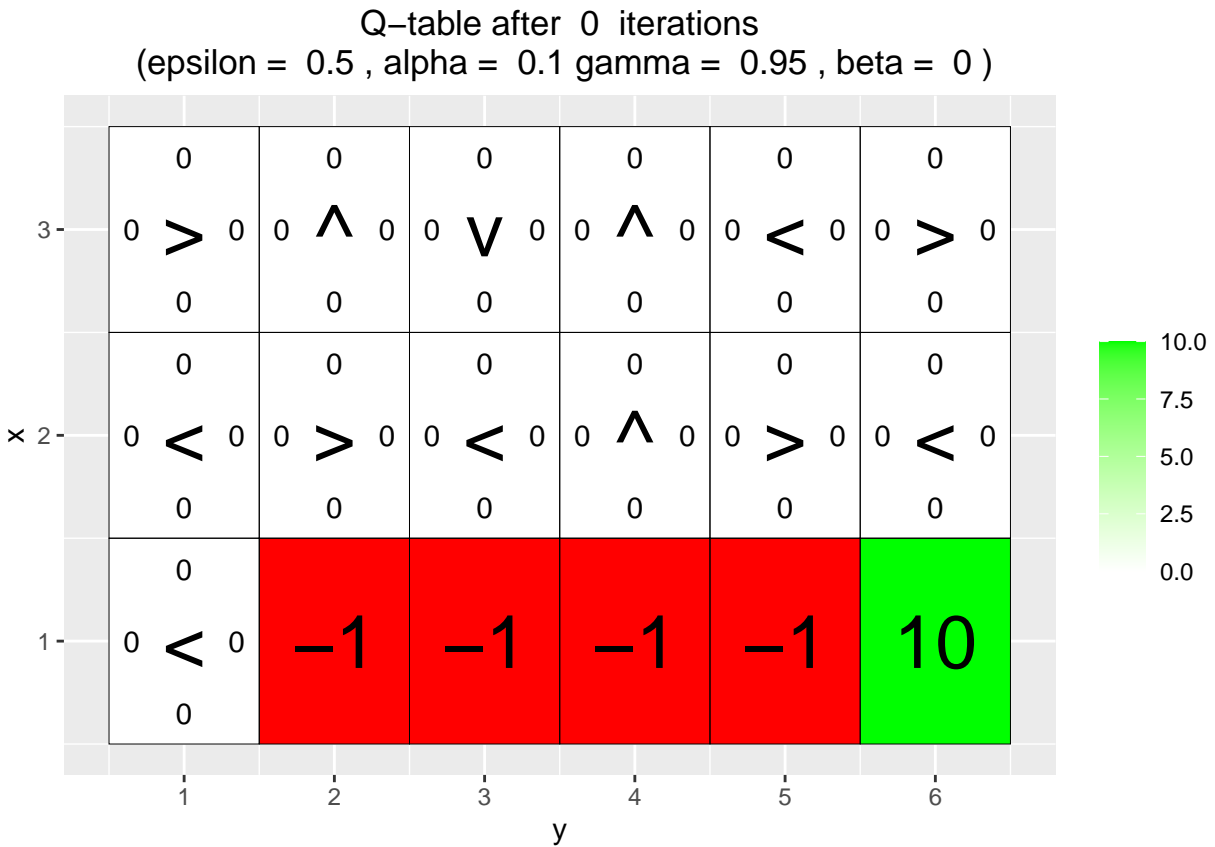
```
# Environment C (the effect of beta).
set.seed(1234)

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()
```

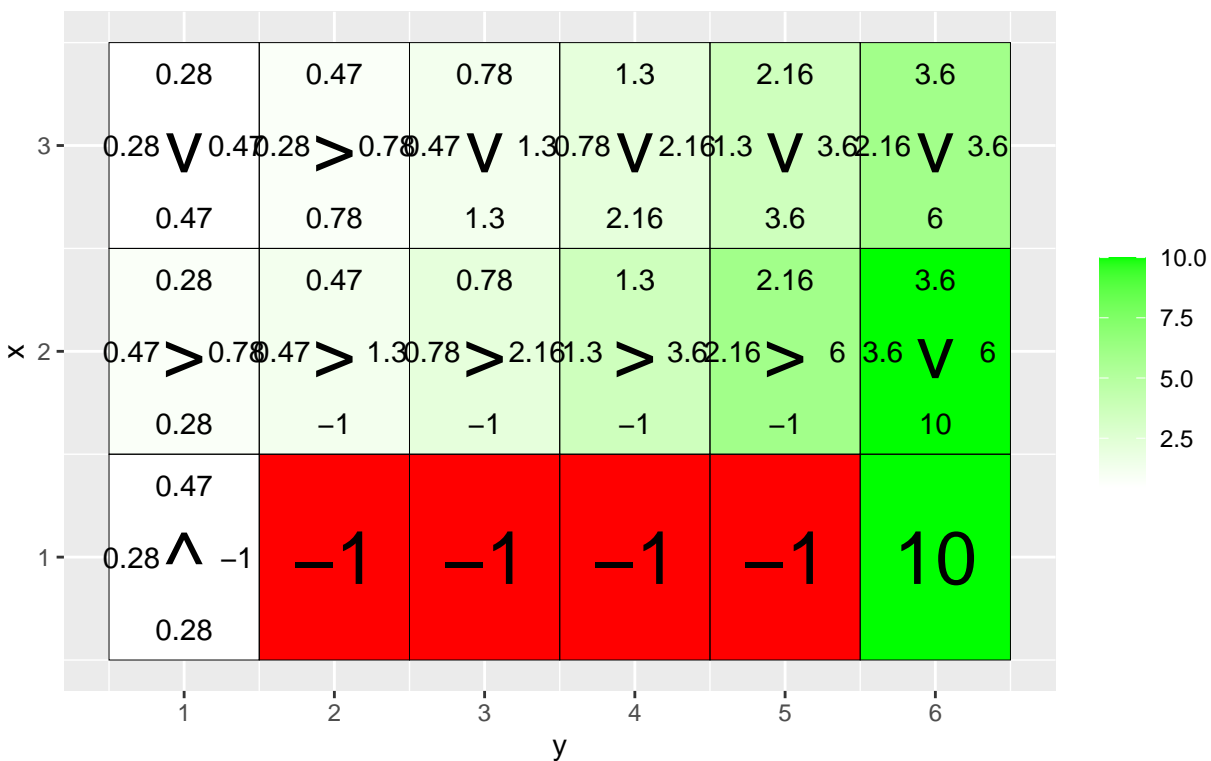


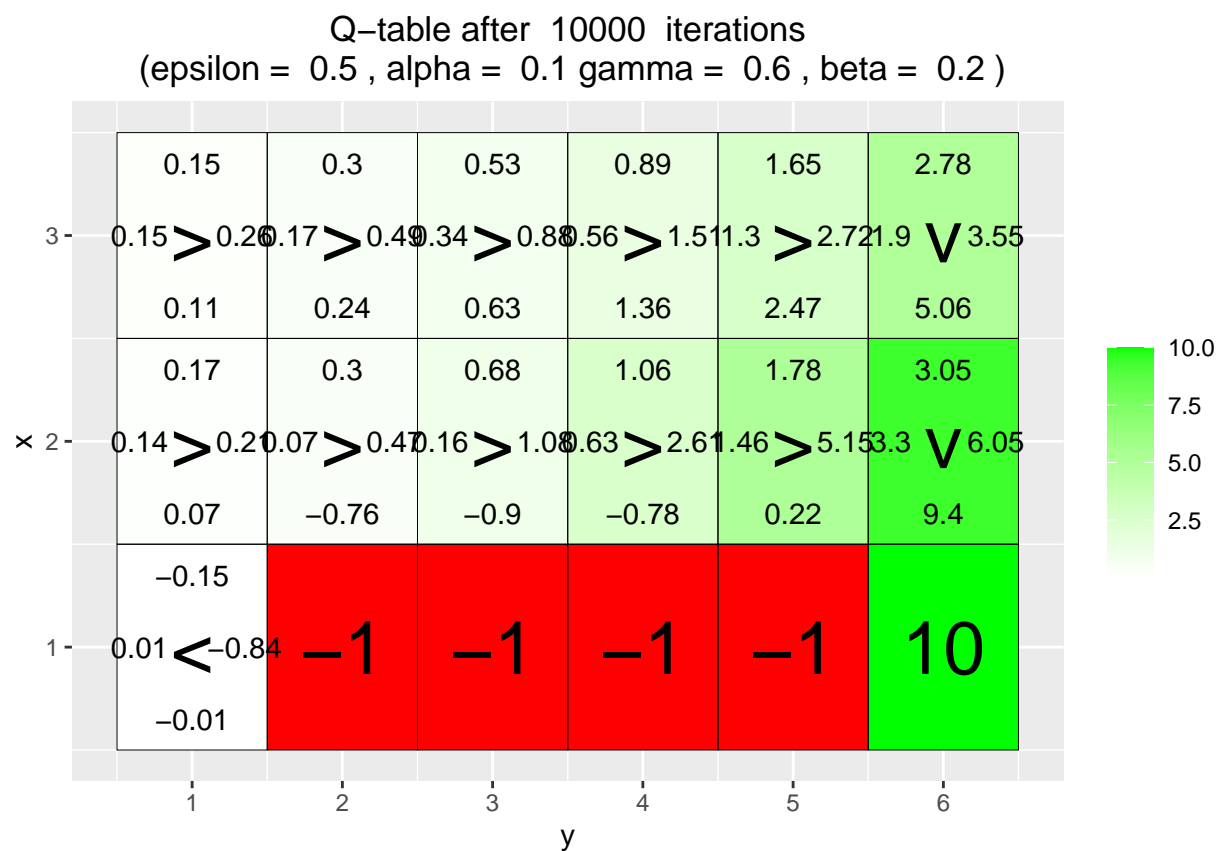
```
for(j in c(0,0.2,0.4,0.66)){
  q_table <- array(0,dim = c(H,W,4))

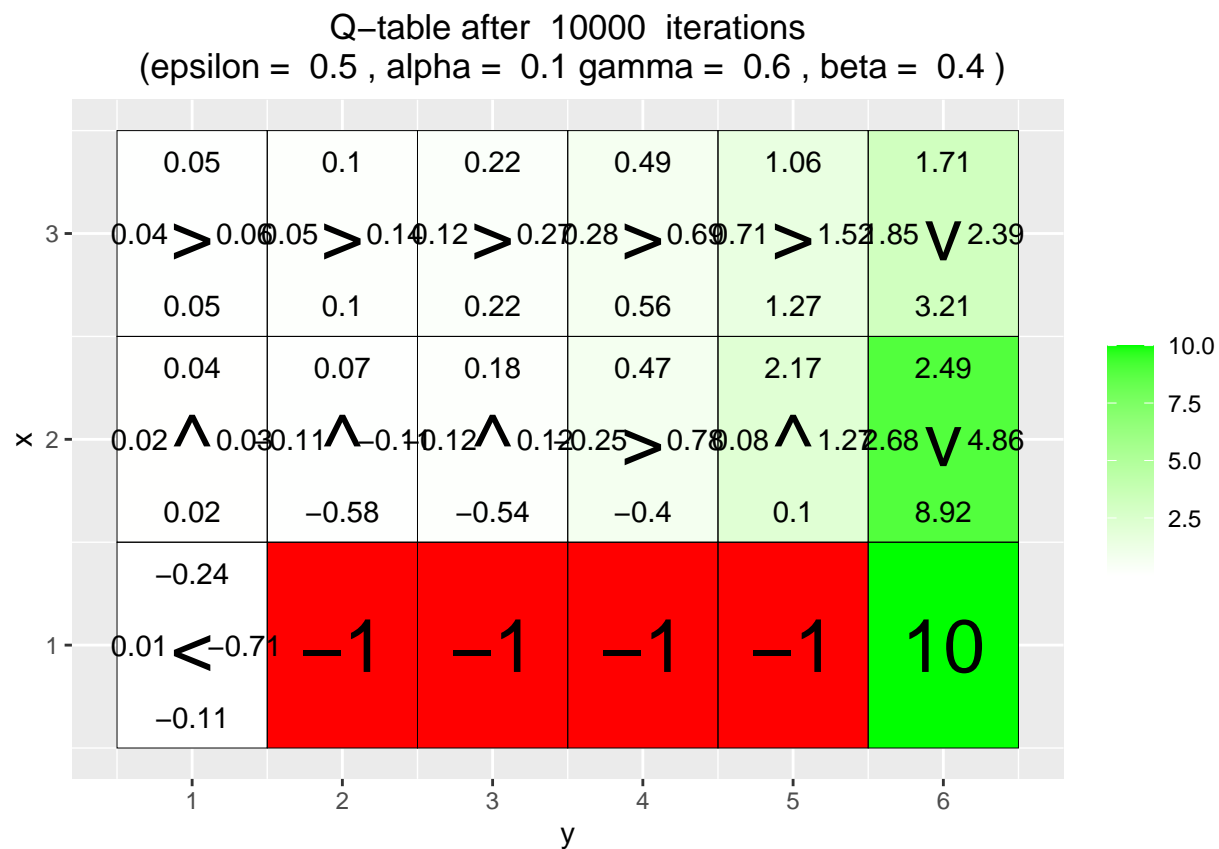
  for(i in 1:10000)
    foo <- q_learning(gamma = 0.6, beta = j, start_state = c(1,1))

  vis_environment(i, gamma = 0.6, beta = j)
}
```

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







[illegible]

3.6

Q: Could you have used the Q-learning algorithm to solve this task ? A: In theory yes, if it got to explore all combinations of states and goals, but Q-learning does not generalize.

Q: Has the agent learned a good policy? Why / Why not ? A: No, the agent tends to want to move upward, since that's what had the best return during training.

Q: If the results obtained for environments D and E differ, explain why. A: The agent in E is only trained on goals to the top, causing its policy to perform bad for goals at the bottom.