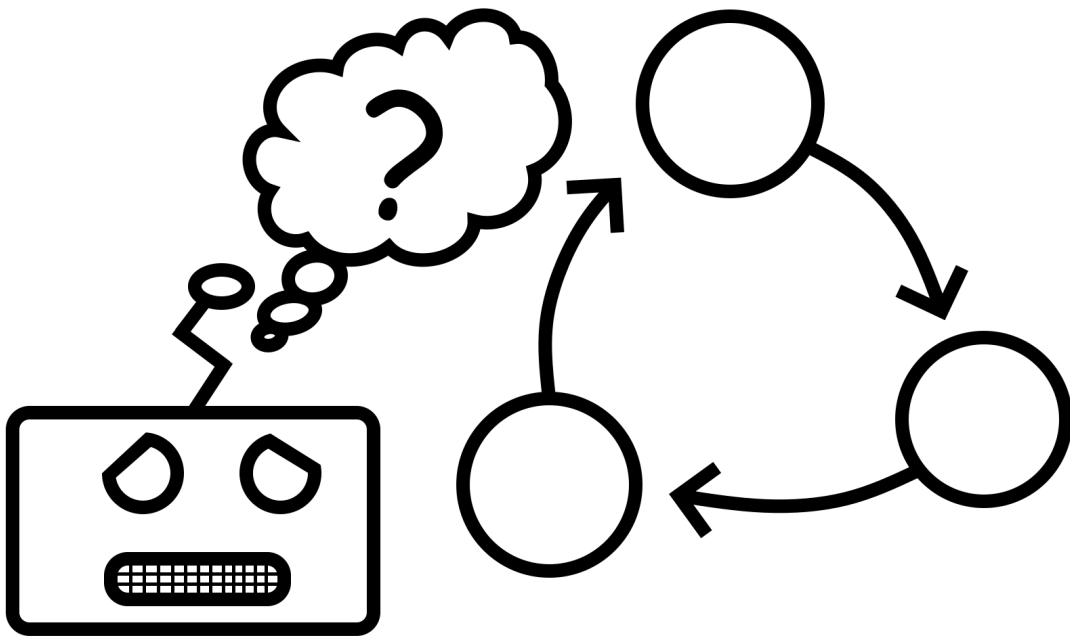


# Model-free Prediction & Control

Programmeeropdracht



John Williams - 1790472

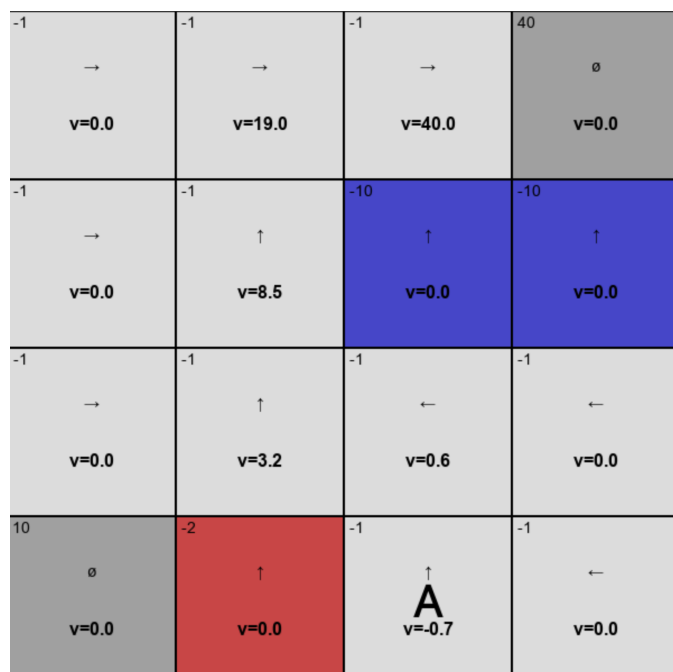
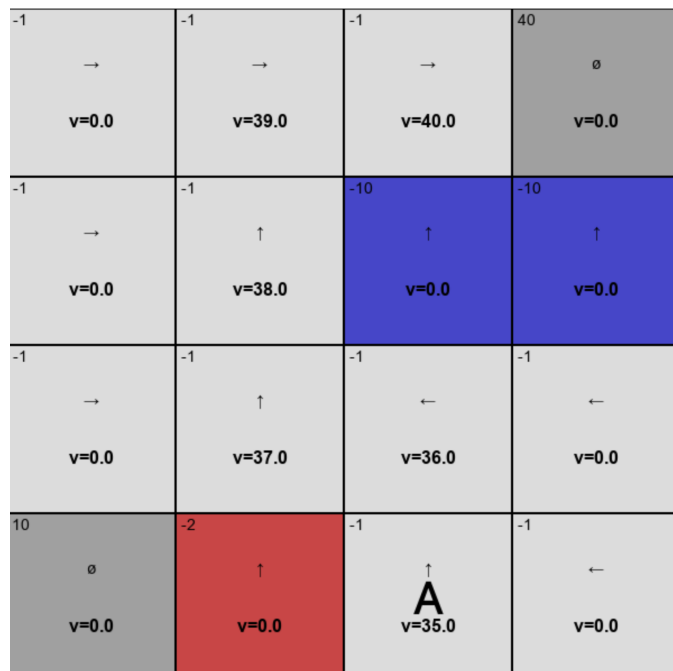
## A. Temporal Difference Learning

1. Implementeer temporal difference learning. Voer de evaluatie uit op de optimale policy  $\pi^*$  met  $\gamma = 1$  en  $\gamma = 0.5$ . Visualiseer de uitkomsten en verklaar het resultaat.

Hier eerst de uitvoering van de evaluatie met  $\gamma = 1$ . Via de pijlen is de optimale policy  $\pi^*$  te zien. Wat opvalt is dat alleen de states direct in het pad van de optimale policy een veranderde waarde hebben (tegenover de initialisatiewaarde 0). Deze waarden verschillen ieder met 1 vergeleken met de volgende/vorige state die door de policy bezocht wordt. De state naast de terminal rechtsboven verwacht de volle reward van de terminal state, terwijl iedere state daarvoor steeds 1 minder in hun reward verwachten, omdat de states zelf een reward van -1 hebben.

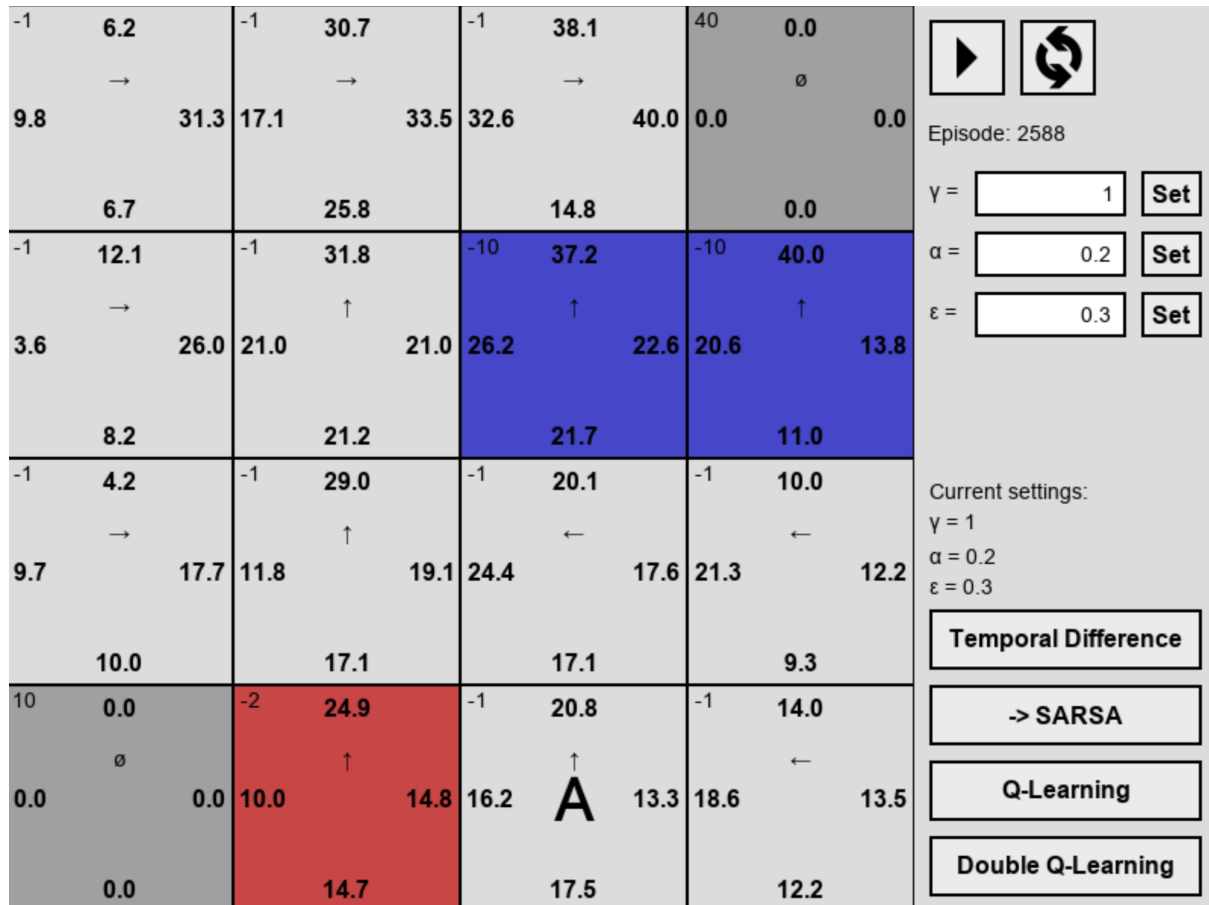
Met een discount van 1 zijn alle toekomstige stappen even belangrijk, dus vandaar dat bij de startpositie al een hoge reward verwacht wordt omdat het een paar stappen van de reward van 40 af is.

Wanneer een discount  $\gamma = 0.5$  wordt gebruikt is het veel duidelijker dat in de toekomst kijken niet veel waarde heeft. De aantallen waarmee de verwachte reward daalt is significant, tot het punt dat de verwachte reward van de startpositie slechter is dan hoe het geïnitieerd was.



## B. SARSA

- I. Implementeer SARSA (on-policy TD control). Voer control uit met  $\gamma = 1$  en  $\gamma = 0.9$ , visualiseer de uitkomsten en verklaar het resultaat.



Hier, waarbij een discount 1 gebruikt is, is te zien dat de policy zich tot de optimale policy vormt. Er is gekozen voor een learning rate van 0.2, om beter de optimale waarden te kunnen benaderen. Er is ook gekozen voor een exploration rate van 0.3. Wat opvalt is dat de agent aan het begin heel graag via de rode state naar de terminal state linksonder gaat. Dit is namelijk een pad die heel makkelijk te bereiken is vanaf het beginpunt. Echter na vele episodes wordt de agent toch consistent in het bereiken van de terminal state rechtsboven.

-1 10.1 → 13.1 23.4 10.2	-1 18.1 → 12.4 29.8 16.1	-1 33.1 → 24.8 40.0 15.0	40 0.0 ∅ 0.0 0.0
-1 7.4 → 3.0 17.0 4.5	-1 23.1 ↑ 9.9 11.1 12.5	-10 33.0 ↑ 17.6 14.6 10.7	-10 40.0 ↑ 14.6 19.4 12.4
-1 5.9 ↓ 5.1 7.0 10.0	-1 17.5 ↑ 7.1 11.0 8.1	-1 9.4 → 12.3 15.1 7.6	-1 20.9 ↑ 10.3 13.3 10.2
10 0.0 ∅ 0.0 0.0	-2 13.5 ↑ 10.0 7.9 7.9	-1 11.2 ↑ 7.6 A 8.6	-1 15.5 ↑ 8.0 9.0 9.4

▶

↺

Episode: 5033

$\gamma =$   Set

$\alpha =$   Set

$\epsilon =$   Set

Current settings:  
 $\gamma = 0.9$   
 $\alpha = 0.2$   
 $\epsilon = 0.3$

Temporal Difference

→ SARSA

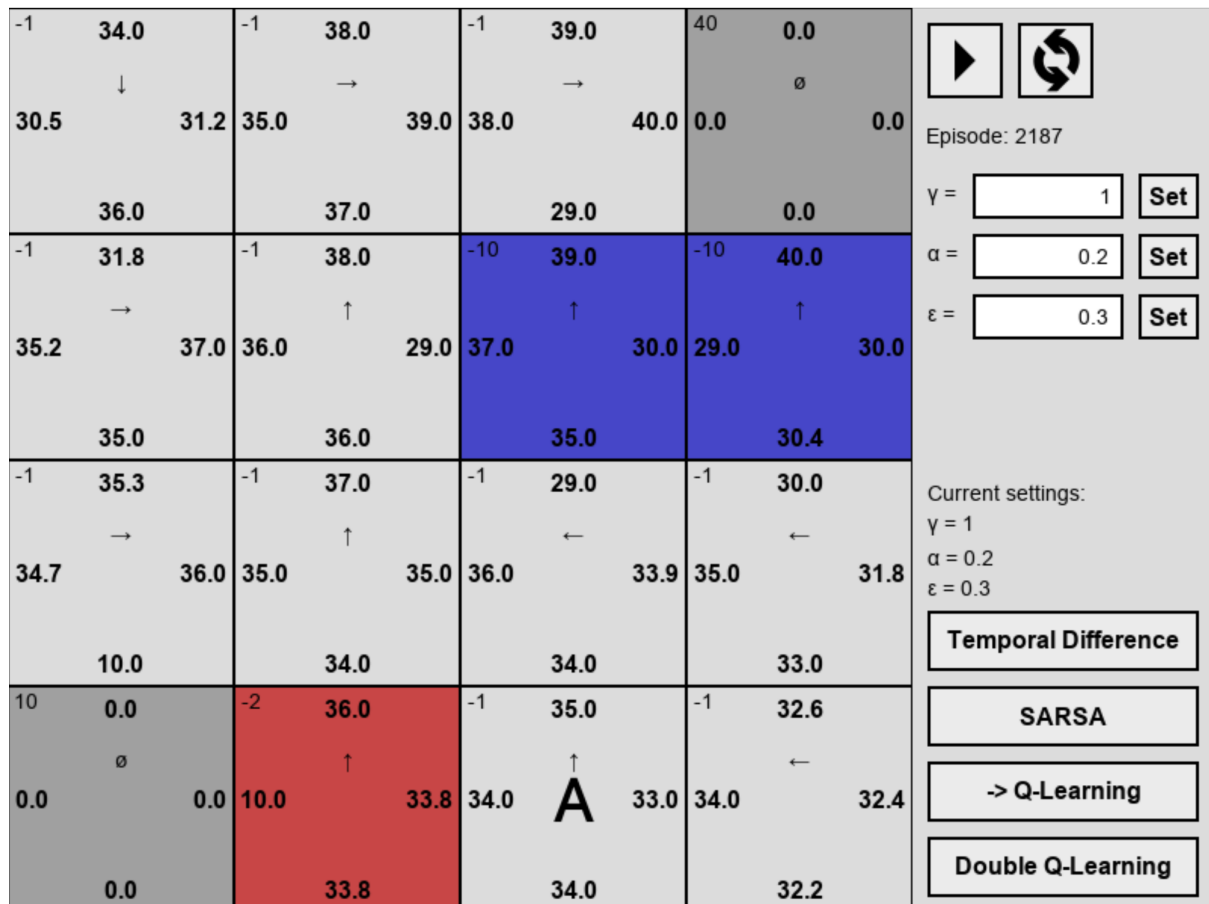
Q-Learning

Double Q-Learning

Hier wordt een discount gebruikt van 0.9, en na dubbel zo veel episodes als de situatie hierboven is het nog steeds niet consistent in het volgen van de optimal policy. De agent heeft de neiging om toch door de blauwe states te gaan, ten kostte van een hogere reward.

## C. Q-Learning

- I. Implementeer Q-learning (off-policy TD control). Voer control uit met  $\gamma = 1$  en  $\gamma = 0.9$ , visualiseer de uitkomsten en verklaar het resultaat.



Hier is Q-Learning uitgevoerd met een discount van 1, en wat meteen opvalt is dat de values van de acties die de optimal policy geconvergeerd zijn tot die van de optimal policy zelf. Dit verschilt met SARSA omdat SARSA een actie kiest van de volgende state op een  $\epsilon$ -greedy manier, terwijl Q-Learning altijd de beste actie van de volgende state kiest. Dit betekent dat wanneer het de optimale actie kiest, dat die waarde geüpdate wordt met behulp van de optimale actie in de volgende state. SARSA update deze waarde met mogelijk de optimale, maar ook mogelijk een suboptimale. Hierdoor heeft SARSA meer moeite met convergeren dan Q-Learning.

-1 0.4 ↓ -0.6 -0.6 9.0	-1 22.1 → 3.2 35.0 15.1	-1 35.0 → 30.5 40.0 21.5	40 0.0 ∅ 0.0 0.0
-1 1.3 → 6.4 16.1 8.0	-1 19.9 → 7.8 21.5 15.8	-10 35.0 ↑ 18.3 26.0 19.1	-10 40.0 ↑ 21.5 26.0 22.4
-1 6.7 ↓ 7.7 5.6 10.0	-1 18.2 → 8.0 19.2 12.2	-1 21.5 → 16.2 22.4 16.2	-1 26.0 ↑ 19.2 22.4 19.2
10 0.0 ∅ 0.0 0.0	-2 15.8 → 10.0 16.2 12.6	-1 19.2 ↑ 12.6 19.2 16.2	-1 22.4 ↑ 16.2 19.0 19.1

Episode: 1756  
 γ =  Set  
 α =  Set  
 ε =  Set

Current settings:  
 γ = 0.9  
 α = 0.2  
 ε = 0.3

Temporal Difference

SARSA

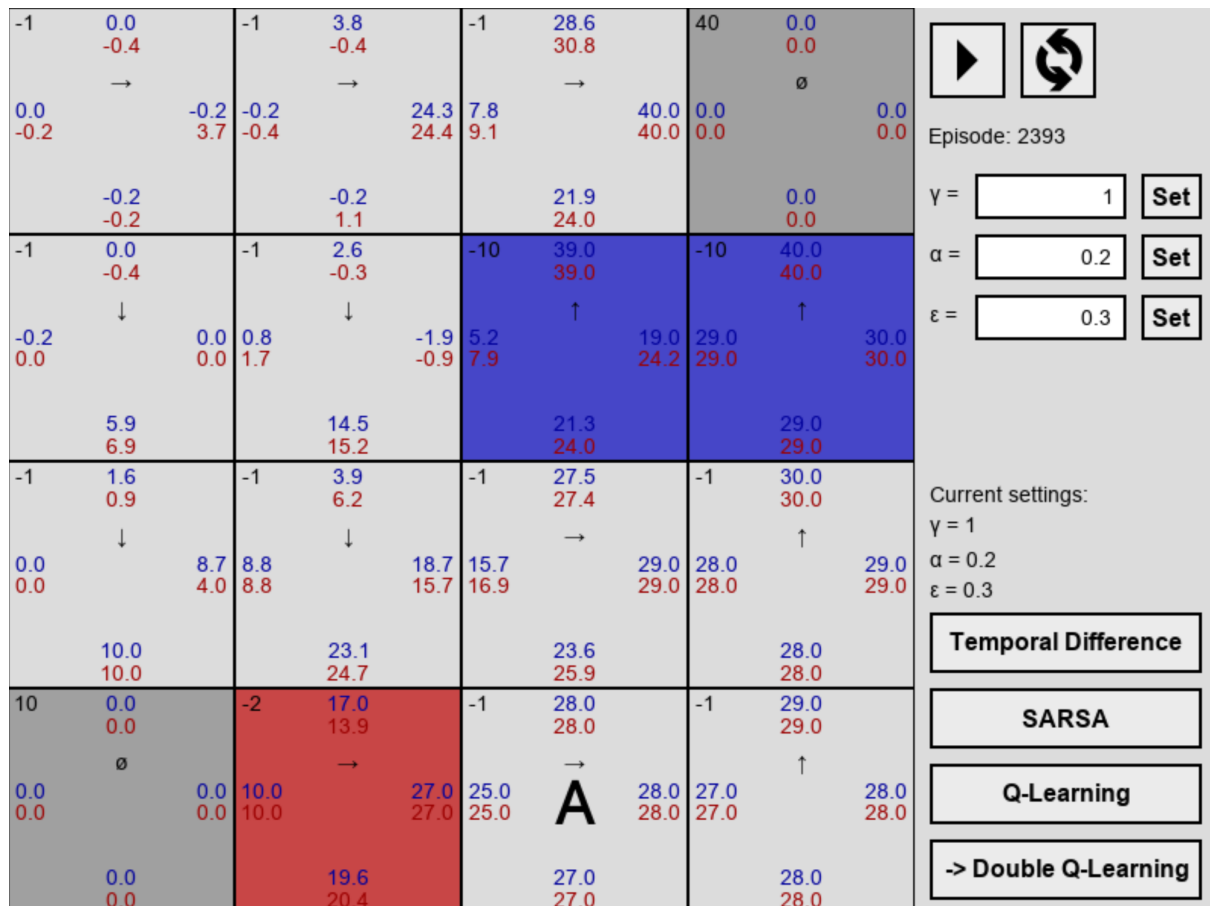
-> Q-Learning

Double Q-Learning

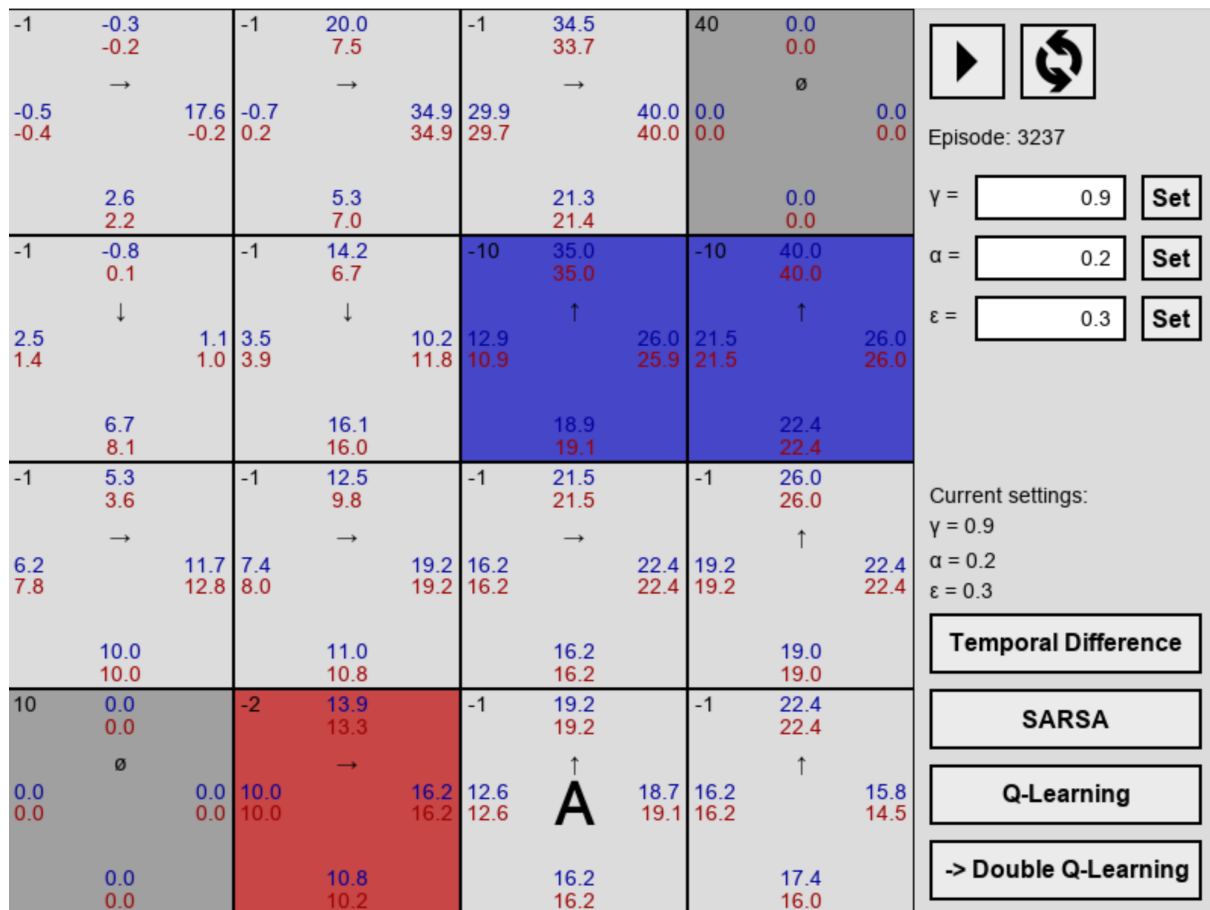
Hier, met een discount van 0.9, is de policy geconvergeerd tot een andere dan de optimal, namelijk eentje die de voorkeur heeft om via de blauwe states te gaan. De waardes veranderen ook niet meer, dus dit is met een discounted vooruitzicht het beste wat je gaat krijgen.

# E. Double Q-Learning

- I. Implementeer Double Q-learning. Voer control uit met  $\gamma = 1$  en  $\gamma = 0.9$ , visualiseer de uitkomsten en verklaar het resultaat. Vergelijk de resultaten tenslotte met de resultaten van Q-learning.



Hier, met een discount van 1, wordt Double Q-Learning uitgevoerd. Wat meteen opvalt vergeleken met de reguliere Q-Learning met een discount van 1, is dat Double Q-Learning niet heeft geconvergeerd tot de optimale policy. Hoe dit gebeurt is kan ik mijzelf momenteel niet verklaren.



Hier is Double Q-Learning uitgevoerd met een discount van 0.9, en net als bij een discount van 0.9 is de policy niet geconvergeerd tot de optimale policy. Deze lijkt echter wel op zijn gelijke bij de reguliere Q-Learning.