



Opdracht 2.2

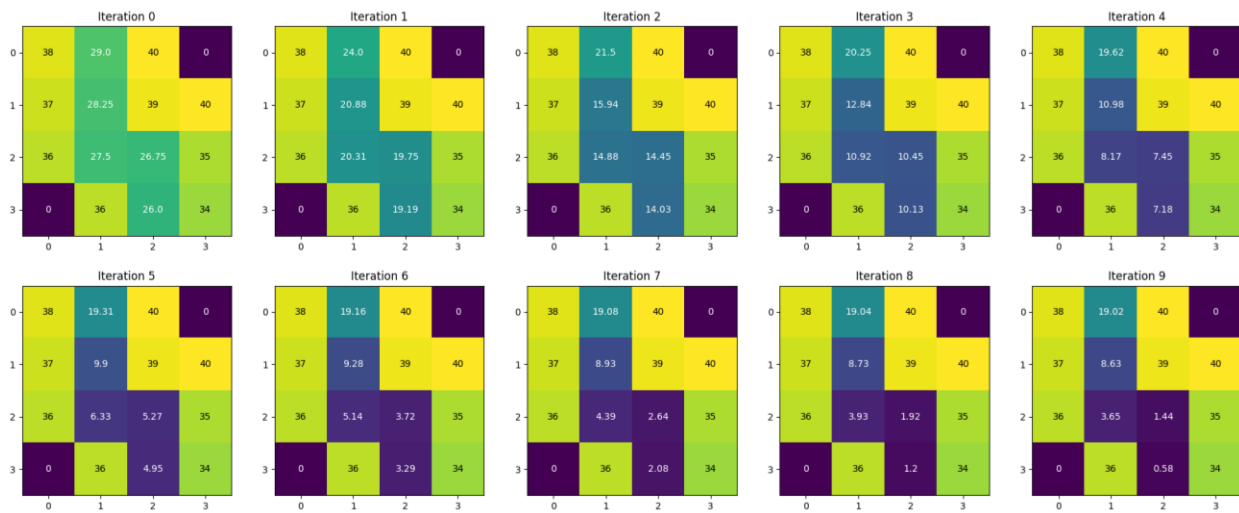
Adaptive Systems

Student:	Storm Joannes
Studentnummer:	1760581
Opleiding:	HBO-ICT Artificial Intelligence
Instelling:	Hogeschool Utrecht
Code:	2022_TICT_VINNO1-33_3_V
Datum:	09-04-2024

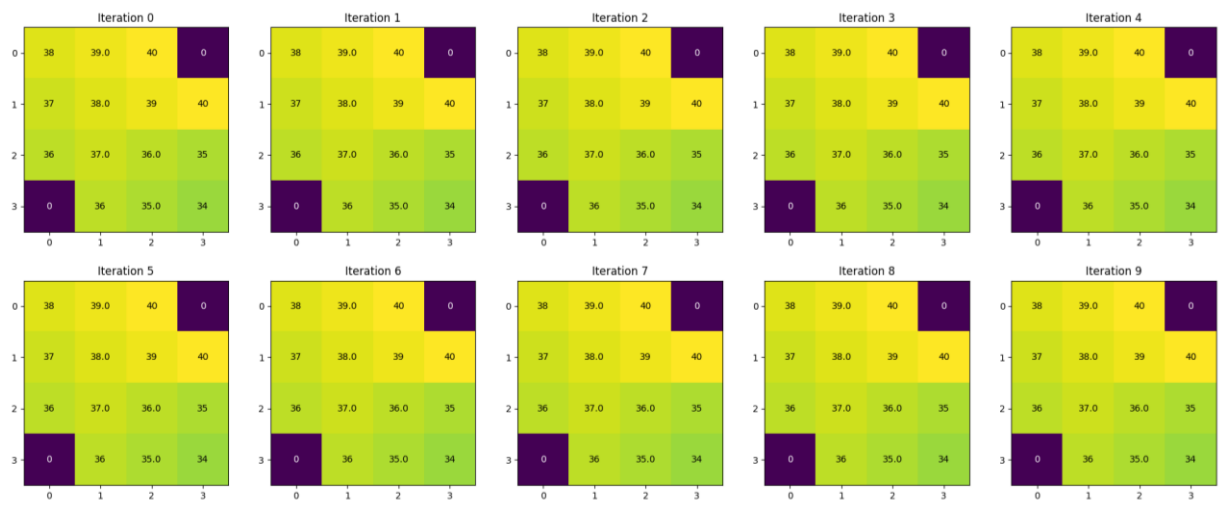
A. Temporal difference learning

De temporal differences worden berekend met 10 epochs, en een learning rate van 0,5.

De eerste figuur zijn de eind waarden van iedere episode temporal difference learning met een discount van 0.5. De episode werd beëindigd als de current position een terminal state was. Dit was in het geval van iedere iteratie positie (0, 3).



In de volgende figuur zijn de iteraties te zien van Temporal difference learning met een discount van 1. Ook hier eindigde iedere iteratie bij de terminal state (0, 3).



Het verschil tussen deze twee grafieken is duidelijk te zien door het verschil in waarden van de afgelopen posities. Het is duidelijk te zien dat de agent met de hogere discount veel sneller convergeert. Dit is te zien doordat de waarden in de plots vanaf de eerste iteratie al niet meer veranderen. In tegenstelling tot de agent met de lagere discount, neemt deze veel langer de tijd om te convergeren, en zoekt deze meer balans tussen de korte- en langetermijnbeloningen.

B. SARSA

De instellingen waarmee wij dit model runnen is een learning rate van 0.25, een epsilon van 0.1 en het aantal epochs van 20.

De eerste is SARSA met een discount van 1.0:

Surrounding values after SARSA_1

(0, 0)	28.9	29.68	29.03	35.4
(0, 1)	36.32	31.56	31.57	38.98
(0, 2)	36.9	26.63	32.81	40.0
(0, 3)	0.0	0.0	0.0	0.0
(1, 0)	30.76	19.23	28.8	34.79
(1, 1)	37.9	30.91	31.55	28.72
(1, 2)	38.63	26.39	31.85	27.97
(1, 3)	40.0	24.31	26.25	21.62
(2, 0)	28.97	10.0	28.09	33.05
(2, 1)	36.7	31.02	29.23	29.7
(2, 2)	27.73	25.13	35.36	27.72
(2, 3)	26.4	24.61	32.58	27.11
(3, 0)	0.0	0.0	0.0	0.0
(3, 1)	34.57	24.83	10.0	24.51
(3, 2)	33.95	29.24	30.52	25.42
(3, 3)	23.26	23.64	29.82	22.49

Up Down Left Right
Next position direction value

Highest Value Direction for Each Coordinate

0	Right	Right	Right	terminal
1	Right	Up	Up	Up
2	Right	Up	Left	Left
3	terminal	Up	Up	Left

Y-axis position
X-axis position

De tweede is SARSA met een discount 0.9:

Surrounding values after SARSA_0.9

(0, 0)	17.89	18.27	16.95	25.02
(0, 1)	19.78	22.51	16.24	32.81
(0, 2)	34.94	16.69	18.14	40.0
(0, 3)	0.0	0.0	0.0	0.0
(1, 0)	23.02	15.26	14.68	13.74
(1, 1)	22.79	16.36	16.84	18.24
(1, 2)	31.59	15.61	19.43	19.2
(1, 3)	40.0	14.01	20.28	26.0
(2, 0)	17.45	10.0	13.39	13.66
(2, 1)	22.88	15.12	14.83	14.03
(2, 2)	11.29	11.96	17.21	14.2
(2, 3)	13.2	9.54	14.24	12.81
(3, 0)	0.0	0.0	0.0	0.0
(3, 1)	18.69	13.88	10.0	10.99
(3, 2)	15.22	13.87	15.58	12.28
(3, 3)	12.76	11.01	12.17	11.16

Up Down Left Right
Next position direction value

Highest Value Direction for Each Coordinate

0	Right	Right	Right	terminal
1	Up	Up	Up	Up
2	Up	Up	Left	Left
3	terminal	Up	Left	Up

Y-axis position
X-axis position

Als we kijken naar de verschillen tussen de twee uitkomsten, zien we bij gekozen acties niet heel veel verschil. We zien dat beide een optimale policy zijn waarbij de agent vanuit de startpositie op een optimale manier je terminal state op $(0, 3)$ weet te bereiken.

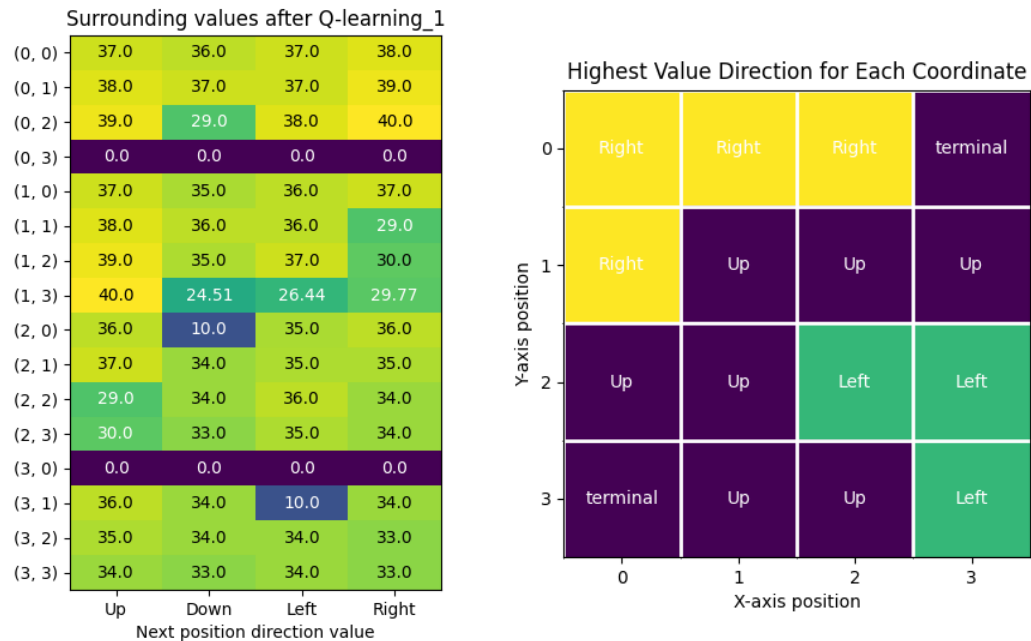
Echter zien we wel een groot verschil als we kijken naar alle waardes van de states. We zien dat met een discount van 0.9 de waardes veel lager zijn. Dit kan komen omdat het bij een discount van 0.9 de agent toekomstige beloningen minder waardeert dan bij een hogere discount. Hierdoor zullen waardes ook lager uitslaan.

De resultaten die hieruit zijn gekomen zijn naar verwachting, en laten geen aparte uitkomsten zien.

C. Q-learning

We hebben hier Q-learning ook wel bekend als SARSA MAX toegepast met een learning rate van 0.25, en een epsilon van 0.1 en het aantal epochs van 30.000, na dit aantal epochs veranderde de resultaten niet meer. Daarnaast waren de resultaten van een discount van 0.9 en 1 hierbij hetzelfde. We hebben voor iedere positie de waardes van de omliggende posities bekeken waarna we een bepaling konden maken van welke stappen de agent zou nemen in de matrix.

Uit Q-learning met een discount van 1:



We kunnen hier zien dat de posities weergegeven in de plot acties weergegeven om zich naar de terminal state te begeven. Iedere positie leidt tot de optimale policy.

Uit Q-learning met een discount van 0.9:

Surrounding values after Q-learning_0.9

(0, 0)	26.44	22.8	26.44	30.5
(0, 1)	30.5	26.45	26.45	35.0
(0, 2)	35.0	21.5	30.5	40.0
(0, 3)	0.0	0.0	0.0	0.0
(1, 0)	26.45	19.52	22.8	26.45
(1, 1)	30.5	22.8	22.8	21.5
(1, 2)	35.0	19.52	26.45	26.0
(1, 3)	40.0	22.4	21.5	26.0
(2, 0)	22.8	10.0	19.52	22.8
(2, 1)	26.45	18.52	19.52	19.52
(2, 2)	21.5	16.57	22.8	22.4
(2, 3)	26.0	19.16	19.52	22.4
(3, 0)	0.0	0.0	0.0	0.0
(3, 1)	22.8	18.52	10.0	16.57
(3, 2)	19.52	16.57	18.52	19.16
(3, 3)	22.4	18.97	16.57	19.15
	Up	Down	Left	Right

Next position direction value

Highest Value Direction for Each Coordinate

0	Right	Right	Right	terminal
1	Right	Up	Up	Up
2	Right	Up	Left	Up
3	terminal	Up	Up	Up
	0	1	2	3

X-axis position

Ook hier zien we weer dat er een optimale policy is uitgekomen. Daarnaast zien we netals bij SARSA dat door de lagere discount de values van de actie states ook lager zijn. Toch kom je wel weer vanaf je start positie (3, 2) op de optimale manier naar de terminal state op (0, 3)