

## Statistical Learning Engineers Final Project

### Muchen Li and Julian Brown

This report contains three tasks. In Task 1 three different algorithms are applied to either regression or classification problem. Task 2 and Task 3 pertain to reinforcement learning algorithms.

#### Task 1 (done by Julian)

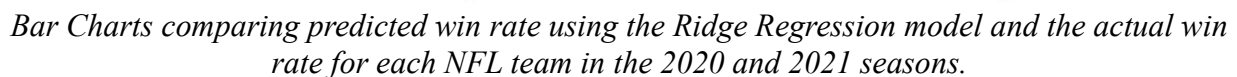
For the Machine Learning Task, a dataset including the 2020 offensive statistics of all 32 teams in the NFL was analyzed along with their win rate over the course of that season, in order to attempt to predict each teams win rate over the current season, 2021. Below are two tables with the data that were used.

Rk	Tm	Tot Yds & TD										Passing					Rushing					Penalties					
		G	PF	Yds	Ply	Y/P	TD	FL	1stD	Cmp	Att	Yds	TD	Int	NY/A	1stD	Att	Yds	TD	Y/A	1stD	Pen	Yds	1stPy	Sc%	TO%	EXP
1	Green Bay Packers	16	31.8	389.0	61.9	6.3	0.69	0.38	22.4	23.3	32.9	256.6	3.00	0.31	7.5	13.5	27.7	122.4	1.00	4.8	7.13	5.23	42.8	1.75	49.7	5.6	18.2
2	Buffalo Bills	16	31.3	396.4	64.6	6.1	1.38	0.69	24.8	25.6	37.3	288.8	2.50	0.69	7.4	15.0	25.7	107.7	1.00	4.2	7.44	6.38	58.8	2.38	48.4	11.8	14.3
3	Tampa Bay Buccaneers	16	30.8	384.1	63.6	6.0	1.06	0.31	22.8	25.6	39.1	289.1	2.63	0.75	7.1	14.9	23.1	94.9	1.00	4.1	5.13	5.23	44.7	2.75	47.8	8.9	15.4
4	Tennessee Titans	16	30.7	396.4	64.4	6.2	0.75	0.31	23.8	19.8	30.3	228.3	2.06	0.44	7.2	12.7	32.6	168.1	1.63	5.2	8.88	5.38	48.9	2.25	47.9	7.2	15.3
5	New Orleans Saints	16	30.1	376.4	63.3	5.8	1.06	0.56	22.9	23.1	32.6	234.9	1.75	0.50	6.8	12.4	30.9	141.6	1.88	4.6	9.19	6.13	62.8	1.31	45.5	9.0	11.1
6	Kansas City Chiefs	16	29.6	415.8	66.1	6.3	1.00	0.56	24.8	26.3	39.4	303.4	2.50	0.44	7.4	15.9	25.2	112.4	0.81	4.5	6.88	6.56	57.4	2.00	47.9	9.2	17.6
7	Baltimore Ravens	16	29.3	363.1	62.1	5.9	1.13	0.44	20.4	16.1	25.4	173.2	1.69	0.65	6.3	8.88	24.7	191.9	1.50	5.5	10.4	6.50	60.1	1.19	45.3	10.6	9.27
8	Seattle Seahawks	16	28.7	369.5	63.9	5.8	1.13	0.31	22.3	24.3	35.2	245.3	2.50	0.51	6.5	13.5	25.7	123.2	0.94	4.8	6.94	5.25	41.4	1.81	44.4	10.1	8.08
9	Indianapolis Colts	16	28.2	378.1	64.5	5.9	0.94	0.25	22.8	23.2	34.5	253.3	1.50	0.69	7.1	12.6	28.7	124.8	1.25	4.3	8.06	5.88	56.2	2.13	44.4	8.8	10.5
10	Las Vegas Raiders	16	27.1	383.3	64.8	5.9	1.63	1.00	22.4	23.1	34.4	263.6	1.75	0.63	7.3	12.9	28.6	119.8	1.25	4.2	7.56	6.13	53.5	2.00	49.9	14.8	9.23
11	Minnesota Vikings	16	26.9	393.3	63.9	6.2	1.44	0.63	23.9	21.8	32.3	250.6	2.19	0.81	7.2	13.3	29.3	142.7	1.25	4.9	8.69	5.13	40.6	2.00	39.8	12.5	9.86
12	Pittsburgh Steelers	16	26.0	334.6	63.2	5.1	1.13	0.44	20.1	26.8	41.0	250.2	2.19	0.69	6.0	12.9	23.3	84.4	0.75	3.6	5.06	5.06	43.2	2.19	37.2	8.5	4.89
13	Arizona Cardinals	16	25.6	384.6	67.7	5.7	1.31	0.50	23.8	24.2	35.9	244.8	1.69	0.81	6.5	13.2	29.9	139.8	1.38	4.7	8.50	7.06	54.3	2.13	40.2	11.7	6.48
14	Cleveland Browns	16	25.5	369.6	63.9	5.8	1.00	0.50	22.2	19.7	31.3	221.2	1.69	0.50	6.7	12.2	30.9	148.4	1.31	4.8	8.31	6.25	54.4	1.69	40.6	9.7	9.97
15	Miami Dolphins	16	25.3	339.0	63.8	5.3	1.25	0.44	21.6	23.1	34.9	233.5	1.50	0.81	6.3	12.9	26.8	105.5	0.94	3.9	6.25	4.83	39.7	2.38	41.4	11.6	3.87
16	Atlanta Falcons	16	24.8	368.4	67.4	5.5	1.13	0.44	22.9	25.5	39.3	272.7	1.69	0.69	6.9	15.2	25.6	95.8	0.81	2.7	5.38	5.19	46.0	2.31	44.8	10.3	8.11
17	Dallas Cowboys	16	24.7	371.8	69.6	5.3	1.63	0.81	23.2	25.8	39.9	260.1	1.56	0.81	6.1	14.4	26.9	111.8	0.88	4.2	7.19	6.00	53.1	1.56	40.6	14.4	5.69
18	Houston Texans	16	24.0	375.3	58.8	6.4	1.13	0.69	20.4	23.9	34.1	283.6	2.06	0.44	7.6	13.9	21.5	91.6	0.63	4.3	5.19	5.00	38.7	1.31	42.9	9.8	6.76
19	Los Angeles Chargers	16	24.0	382.1	70.4	5.4	1.00	0.38	23.3	25.8	39.2	270.6	1.94	0.63	6.5	14.1	29.1	111.5	0.75	3.8	6.94	5.31	44.4	2.25	38.3	8.6	8.10
20	Detroit Lions	16	23.6	350.2	61.9	5.7	1.31	0.50	21.9	23.4	36.4	256.5	1.69	0.81	6.6	13.8	22.9	93.7	1.06	4.1	5.81	5.94	53.8	2.31	38.5	11.2	5.92
21	San Francisco 49ers	16	23.5	370.1	65.4	5.7	1.94	0.88	21.9	23.2	35.6	252.1	1.56	1.06	6.6	13.6	27.3	118.1	1.19	4.3	6.31	5.31	45.7	2.00	36.3	14.8	3.48
22	Los Angeles Rams	16	23.3	377.0	68.0	5.5	1.56	0.69	22.0	24.5	36.9	250.9	1.25	0.88	6.5	12.9	29.6	126.1	1.19	4.3	7.69	4.44	40.9	1.44	34.2	12.5	3.58
23	Chicago Bears	16	23.3	331.4	65.2	5.1	1.38	0.38	20.9	25.1	38.4	228.4	1.63	1.00	5.6	13.6	24.6	102.9	0.75	4.2	5.81	5.50	48.6	1.50	38.9	12.0	11.8
24	Carolina Panthers	16	21.9	349.5	62.1	5.6	1.31	0.31	20.9	23.3	34.4	245.0	1.00	1.00	6.8	12.0	25.4	106.5	1.19	4.2	6.88	5.94	47.0	2.06	40.3	13.2	5.08
25	Washington Football Team	16	20.9	317.3	65.7	4.8	1.69	0.69	20.1	24.3	37.6	216.6	1.00	1.00	5.3	11.5	25.0	100.7	1.13	4.0	6.75	5.44	43.6	1.88	33.3	13.1	-0.65
26	Philadelphia Eagles	16	20.9	334.6	64.6	5.0	1.81	0.56	21.0	20.9	37.4	207.9	1.38	1.25	5.0	11.1	25.2	126.7	1.00	5.0	7.13	6.69	53.3	2.81	27.9	15.5	-0.85
27	New England Patriots	16	20.4	327.3	61.3	5.2	1.19	0.31	20.8	17.7	27.5	180.6	0.75	0.88	6.1	9.81	31.4	146.6	1.25	4.7	8.94	3.88	33.4	2.00	36.9	12.1	3.49
28	Denver Broncos	16	20.2	335.6	64.4	5.2	2.00	0.56	19.3	19.8	34.8	215.7	1.31	1.44	5.9	11.0	27.6	119.9	0.81	4.3	5.88	5.06	41.9	2.38	32.0	17.0	-1.36
29	Cincinnati Bengals	16	19.4	319.8	65.0	4.9	1.50	0.81	19.9	23.3	36.3	215.5	1.19	0.69	5.5	12.6	25.7	104.3	0.81	4.1	5.75	5.19	42.2	1.56	33.0	12.8	-0.32
30	Jacksonville Jaguars	16	19.1	326.1	62.3	5.2	1.56	0.56	19.4	24.2	38.5	231.2	1.56	1.00	5.6	12.8	21.1	94.9	0.56	4.5	5.00	6.69	66.9	1.56	30.4	14.0	-0.38
31	New York Giants	16	17.5	299.6	60.4	5.0	1.38	0.69	18.6	20.1	32.3	189.1	0.75	0.69	5.3	11.1	24.9	110.5	0.81	4.4	5.69	5.06	39.6	1.75	33.5	12.6	-1.67
32	New York Jets	16	15.2	279.9	59.3	4.7	1.19	0.31	16.8	18.3	31.2	174.8	1.00	0.88	5.2	9.13	25.4	105.2	0.56	4.1	5.88	6.31	59.5	1.81	26.3	10.9	-3.69

*2020 NFL per-game Offensive Statistics – These data were used as training for each model.*

Rk	Tm	Tot Yds & TD										Passing					Rushing					Penalties					
		G	PF	Yds	Ply	Y/P	TD	FL	1stD	Cmp	Att	Yds	TD	Int	NY/A	1stD	Att	Yds	TD	Y/A	1stD	Pen	Yds	1stPy	Sc%	TO%	EXP
1	Tampa Bay Buccaneers	13	31.5	410.2	67.1	6.1	1.23	0.46	24.3	29.6	43.3	314.2	2.77	0.77	7.0	16.3	22.5	96.0	1.08	4.3	6.54	6.31	37.3	1.69	45.0	10.7	13.4
2	Dallas Cowboys	13	29.2	409.1	68.6	6.0	1.38	0.54	22.9	26.1	38.5	280.6	2.00	0.85	6.9	14.3	28.2	128.5	0.85	4.6	6.46	7.92	70.1	2.15	41.1	11.3	6.37
3	Indianapolis Colts	13	28.5	368.1	63.8	5.8	1.23	0.77	22.3	20.6	32.7	216.4	1.69	0.46	6.3	10.7	29.5	151.7	1.54	5.1	9.54	4.77	42.9	2.08	45.1	9.2	5.96
4	Los Angeles Rams	13	28.2	384.5	62.0	6.2	1.08	0.31	21.1	24.4	36.5	287.3	2.54	0.77	7.5	14.2	23.9	97.2	0.62	4.1	5.46	4.92	42.2	1.38	48.6	9.4	10.1
5	Arizona Cardinals	13	28.2	374.8	65.0	5.8	1.00	0.23	21.6	23.4	32.3	252.2	1.69	0.77	7.3	12.6	30.3	122.5	1.62	4.0	7.31	6.38	56.2	1.69	45.3	8.6	9.29
6	Buffalo Bills	13	27.9	382.9	65.6	5.8	1.38	0.46	22.8	25.4	38.6	261.9	2.15	0.92	6.5	13.9	25.3	121.0	1.00	4.8	7.23	6.85	60.1	1.62	44.8	11.9	7.79
7	Cincinnati Bengals	13	27.2	358.8	62.8	5.6	1.62	0.54	19.8	22.1	32.2	250.1	2.00	1.08	7.1	12.4	27.1	108.7	1.15	4.0	5.85	4.00	34.1	1.62	41.1	12.6	4.57
8	Los Angeles Chargers	13	27.0	385.2	64.5	6.0	1.15	0.31	23.1	26.0	38.8	280.3	2.31	0.85	6.8	14.7	23.6	104.9	0.92	4.4	6.38	7.15	61.1	2.00	44.4	11.1	9.38
9	Kansas City Chiefs	13	27.0	389.6	66.7	5.8	1.77	0.85	24.4	26.2	39.9	277.9	2.08	0.92	6.7	15.5	24.9	111.7	0.92	4.5	6.85	6.69	55.4	2.08	45.4	16.2	9.89
10	New England Patriots	13	26.9	346.5	61.4	5.6	1.23	0.42	20.7	21.5	30.3	223.0	1.38	0.42	6.9	10.9	29.2	123.3	1.15	4.2	7.54	5.62	51.2	2.23	48.2	10.9	9.85
11	Minnesota Vikings	13	26.5	390.5	66.2	5.9	0.85	0.46	21.4	24.7	37.0	267.4	2.08	0.38	7.0	12.8	27.8	123.2	0.69	4.4	6.85	7.31	69.2	1.69	42.0	7.2	6.53
12	Philadelphia Eagles	13	25.3	356.2	61.8	5.7	1.23	0.38	20.9	17.9	26.2	246.5	1.69	0.38	6.2	8.6	17.1	107.6	1.00	3.8	6.17	6.08	54.1	1.69	42.0	10.7	9.85
13	Washington Redskins	13	25.3	364.3	61.8	5.7	1.26	0.69	21.0	20.7	30.7	240.4	1.54	0.38	6.2	11.5	29.1	129.3	1.31	4.3	7.54	6.08	62.3	2.00	40.7	11.4	6.37
14	San Francisco 49ers	13	25.2	361.7	62.7	5.8	0.77	0.28	21.5	23.1	34.7	253.8	2.15	0.38	6.9	13.3	26.0	107.6	0.69	4.1	6.15	4.38	42.7	2.0	42.4	6.1	9.38
15	Tennessee Titans	13	24.9	347.2	67.2	5.6	1.62	0.42	21.8	21.4	32.5	299.4	1.15	1.00	7.0	11.8	37.8	148.4	0.44	4.3	7.92	6.08	56.6	2.23	48.2	13.7	4.89
16	Baltimore Ravens	13	23.4	388.0	70.8	6.5	1.54	0.46	23.8	28.1	35.9	242.2	1.38	0.85	6.1	12.6	31.3	144.8	1.08	4.6	9.46	3.68	56.9	1.77	39.0	10.2	14.4
17	New Orleans Saints	13	23.4	319.4	62.4	5.1	1.23	0.38	19.2	18.1	31.2	197.9	1.92	1.05	6.0	10.5	34.9	121.5	0.92	4.1	7.00	5.85	48.8	1.69	34.2	10.1	14.4
18	Los Vegas Raiders	13	21.8	372.5	62.7	5.9	1.31	0.54	20.4	26.1	38.2	287.8	1.38	0.77	7.1	13.0	22.2	84.6	0.83	5.0	7.80	6.85	71.1	2.31	39.0	11.6	41.3
19	Cleveland Browns	13	21.4	349.2	61.6	5.7	1.08	0.31	20.2	18.9	30.2	205.7	1.08	0.64	6.3	10.5	28.8	143.5	1.31	5.0	5.92	6.85	62.5	1.69	34.6	10.5	41.3
20	Denver Broncos	13	21.2	345.0	62.7	5.5	1.15	0.46	20.6	22.2	33.2	221.9	1.38	0.69	6.2	11.3	26.9	123.1	0.92	4.6	7.62	5.00	41.9	1.69	38.2	10.7	40.0
21	Seattle Seahawks	13	20.9	310.2	55.4	5.6	0.77	0.38	17.2	19.7	29.4	205.8	1.62	0.38	6.4	8.77	32.0	104.4	0.92	4.5	5.46	5.62	43.5	1.92	30.3	16.3	6.08
22	Indianapolis Steelers	13	20.9	329.0	64.3	5.1	1.15	0.54	19.4	25.8	38.9	240.6	1.54	0.62	5.8	12.3	27.7	88.4	0.62	3.7	9.42	6.54	51.3	2.15	35.4	9.5	9.38
23	Washington Football Team	13	20.5	336.2	64.6	5.2	1.62	0.69	21.3	21.9	33.3	213.7	1.46	0.92	6.0	11.0	32.0	122.5	0.69	4.2	7.69	5.77	48.2	1.69	35.8	15.3	18.6
24	Carolina Panthers	13	19.8	310.6	64.8	4.8	1.77	0.46	19.5	20.2	34.9	209.5	0.85	1.31	5.5	10.6	27.6	107.7	1.15	4.0	7.08	7.23	58.6	1.77	31.8	15.2	4.77
25	Hiam Dolphins	13	19.5	309.8	64.9	4.8	1.54	0.77	19.0	25.8	38.5	230.5	1.31	0.77	5.6	12.9	32.9	79.2	0.69	3.3	4.42	5.86	48.2	1.62	30.3	12.0	4.77
26	Hiam Packers	13	18.6	316.4	61.9	5.1	1.54	0.46	20.1	23.6	35.5	225.3	1.31	0.80	6.0	12.3	24.5	91.0	0.62	3.7	4.46	5.92	48.9	2.41	37.1	14.4	0.57
27	New York Giants	13	18.2	312.2	62.1	5.1	1.31	0.46	18.9	22.9	36.2	217.9	1.05	1.18	5.7	11.3	23.3	94.2	0.62	4.0	5.38	5.98	42.5	2.23	32.2	11.9	4.77
28	San Diego Chargers	13	17.8	309.5	61.9	5.0	1.23	0.38	18.6	20.7	33.8	217.6	1.69	0.38	6.0	10.5	23.3	104.4	0.92	4.5	5.46	5.62	43.5	1.92	30.3	16.3	6.08
29	Denver Broncos	13	17.4	221.5	62.8	5.1	1.92	0.38	19.2	23.4	38.8	227.4	1.31	0.54	5.7	12.5	21.3	84.2	0.69	3.9	4.69	4.69	54.2	2.00	29.0	17.2	0.57
30	Cincinnati Bengals	13	16.4	311.3	62.0	5.0	1.38	0.42	19.2	23.5	35.6	200.8	0.77	0.77	5.3	11.1	24.0	116.5	0.69	4.4	6.23	6.69	51.8	2.8	29.2	31.3	2.41
31	Indianapolis Jaguars	13	13.8	303.1	60.1	5.0	1.92	0.85	21.4	21.0	34.0	204.0	0.69	1.08	5.3	10.1	22.2	102.7	0.62	4.4	5.38	6.85	56.3	1.62	21.7	37.4	9.38
32	Houston Texans	13	12.4	264.2	59.2	4.5	1.54	0.54	15.9	20.8	32.6	186.8	1.00	1.00	5.3	8.46	22.8	77.5	0.46	3.3	4.20	6.52	56.7	2.41	21.1	12.9	7.79

Comparing the scores of each of the three best models, the Ridge Regression was the best performing model for this application. The plot below shows the predicted and actual win rate for each NFL in the 2021 season as determined by the Ridge Regression.



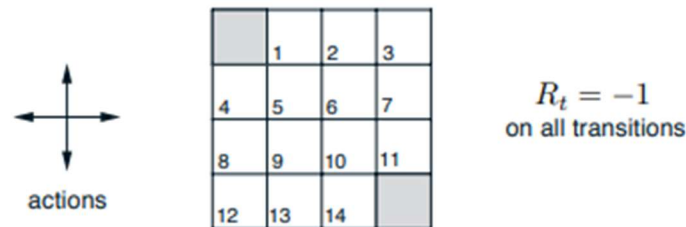
As expected, only using offensive statistics and one season of training data did not result in the most accurate of models for predicting win rate. Offense is, after all, only half of the game (less than half if you include Special Teams play). It is, however, interesting to see the predicted win rate was almost always lower than the actual win rate for the 2021 season. Seeing as the training was done on the 2020 data and that performance was much more evenly balanced in predicted and actual, this indicates that teams in 2021 are winning more with less offensive

production than in 2020. This could mean that better defense and special teams are being played this year, although a number of other factors could also be involved.

Data was accessed at <https://www.pro-football-reference.com/>. Sklearn manuals were accessed online at <https://scikit-learn.org/stable/index.html> for this section.

### Problem for DP (done by Julian)

In this task, we implemented a Policy Iteration loop for Example 4.1 in the RL textbook.



*A screenshot from the RL textbook displaying the actions, states, and rewards.*

As shown above this problem included two terminal states, 14 non-terminal states, 4 equiprobable actions (right, left, up and down), and a reward of  $-1$  for all transitions. The resultant Final Value Function and Optimal Policy, after completing the Policy Iteration algorithm, are shown below.

<b>Final Value Function:</b>	<b>Optimal Policy (right, left, up, down):</b>
<code>[[ 0.   -1.   -1.9  -2.71]</code>	<code>[[ '0'  '←'  '←'  '←']</code>
<code>[-1.   -1.9  -2.71  -1.9 ]</code>	<code>['↑'  '←'  '→'  '↓']</code>
<code>[-1.9  -2.71  -1.9   -1.  ]</code>	<code>['↑'  '→'  '→'  '↓']</code>
<code>[-2.71  -1.9  -1.    0.  ]]</code>	<code>['→'  '→'  '→'  '0']]</code>

*The Final Value Function and Optimal Policy after completion of Policy Iteration.*

Both of these make sense in relation to the setup of this problem, indicating that the algorithm was successful. The policy indicates which way the optimal action is at each state. For states with more than one optimal action, the order of actions shown in the top right was used to determine which arrow to display. For example, the best action for state in position (1,1) would be to move either up or left, but because we chose left to come before up, left is displayed.

## Task 2 Reinforcement Learning (done by Muchen)

In this task, we apply reinforcement learning techniques to a 2-D inverted pendulum. Figure 1 shows a schematic configuration. Specifically, we applied policy iteration, value iteration and q-learning. A complete description can be found in FinalProjectInstruction.pdf.

### Policy Iteration and Value Iteration

This section describes the implementation of two dynamic programming algorithms: policy iteration and value iteration.

The boxes system for the state space is constructed by the following quantization thresholds:

$$\begin{aligned} x: & \pm 0.8, \pm 2.4, \pm \inf \text{ } m \\ x': & \pm 0.5, \pm \inf \text{ } m/s \\ \theta: & 0, \pm 6, \pm 12, \pm \inf^\circ \\ \theta': & \pm 10, \pm 30, \pm 50, \pm \inf^\circ/s \end{aligned}$$

Note that our boxes system is different from that in instruction. The angle velocity space is divided into more boxes because, as a result, the DP model learns faster. In other word, the policy converges with less epochs.

### Physical model

For simplicity, let the cart mass be  $M$ , the pole (pendulum) mass be  $m$  with the half-pole length  $l$ . The model is written as the following nonlinear equations

$$\begin{aligned} \ddot{\theta} &= \frac{g \sin \theta + \cos \theta \left[ \frac{-F - ml\dot{\theta}^2 \sin \theta + \mu_c \text{sgn}(\dot{x})}{M+m} \right] - \frac{\mu_p \dot{\theta}}{ml}}{l \left[ \frac{4}{3} - \frac{m \cos^2 \theta}{m+M} \right]} \\ \ddot{x} &= \frac{F + ml \left[ \dot{\theta}^2 \sin \theta - \ddot{\theta} \cos \theta \right] - \mu_c \text{sgn}(\dot{x})}{M+m} \end{aligned}$$

Figure 1 Non-linear differential equations to model the cartpole system

where  $\text{sgn}$  is the sign function [1]. In implementation, we add coefficients of friction  $\mu_c$  and  $\mu_p$  to the original equations in the CartPoleEnv class as the following [2]:

```

temp = (
    force + self.polemass_length * theta_dot ** 2 * sintheta
) / self.total_mass
temp2 = self.mu_p*theta_dot/self.polemass_length
thetaacc = (self.gravity * sintheta - costheta * (temp - self.mu_c*sgn(x_dot)/self.total_mass) -
temp2) / (
    self.length * (4.0 / 3.0 - self.masspole * costheta ** 2 / self.total_mass)
)
xacc = temp - self.polemass_length * thetaacc * costheta / self.total_mass - self.mu_c*sgn(x_dot) /
self.total_mass

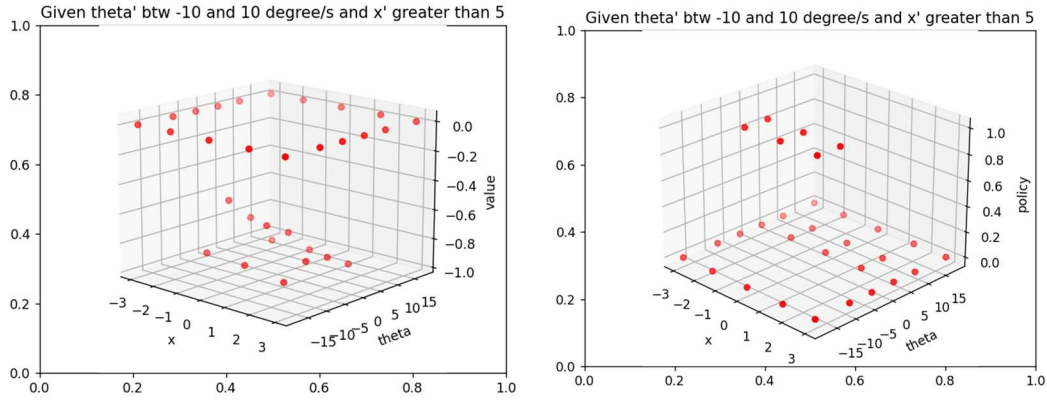
```

The rest of the implementation of the physical model are the same as that in the CartPoleEnv class.

### ***Value Iteration Result***

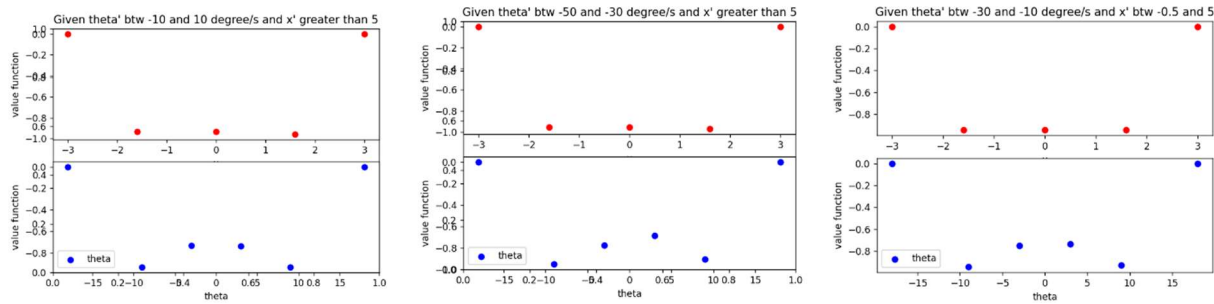
Since the value function is a function of four variables, we fix  $x'$  and  $\theta'$  in a way described in the title of each plot, and then plot value function with respect to  $x$  and  $\theta$ . Plots, policy and value function data are accessible in `plot_task3_value_iteration.py`.

Below, the left figure is value function vs  $x$  and  $\theta$ , and the right figure is policy vs  $x$  and  $\theta$ , given  $\theta' \in (-10, 10)^\circ/s$  and  $x' > 5 \text{ m/s}$ . Noted that the dotted square at value = 0 corresponds to the ‘terminal’ states, i.e., states with pole angle more than 12 degrees or cart position is more than 2.4 meters. The value 0 is the initialized value of the value function at the start of the iteration. Since the only necessary action for the terminal states is ‘exit’, there is no need to update their value function. Also, in the right figure, except the terminal states, when  $\theta < 0$ , the policies are always going right and vice versa. The resulted policy from value iteration helps keep the pole straight and thus is correct.

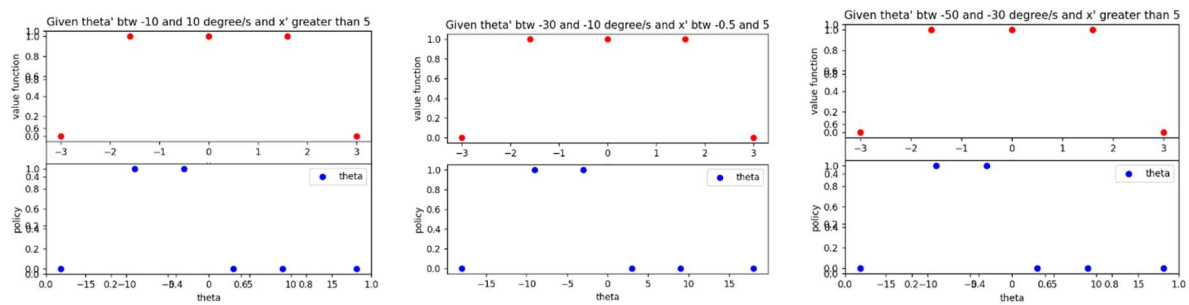


### Policy Iteration Result

Since the value function is a function of four variables, we fix  $x'$  and  $\theta'$  in a way described in the title of each plot, and then fix  $x$  to plot value function with respect to  $\theta$  and vice versa. For example in the left figure, we first select the values corresponding to  $\theta' \in (-10, 10)^\circ/s$  and  $x' > 5 \text{ m/s}$ .



The following are the policy plots. Note that policy 1 is right direction and 0 left.





## Q learning

### The boxes system, learning rate and exploration rate $\epsilon$

The boxes system for the state space is constructed differently than in the DP problem:

$$\begin{aligned} x: & \pm 1.6, \pm \inf m \\ x': & \pm 0.5, \pm \inf m/s \\ \theta: & 0, \pm 8, \pm 16, \pm \inf^\circ \\ \theta': & \pm 10, \pm 30, \pm 50, \pm \inf^\circ/s \end{aligned}$$

The learning rate and the exploration rate decreases logarithmically to 0.1 and 0.001 respectively [3]:

```
def select_learning_rate(x):
    return max(min_learning_rate, min(1.0, 1.0 - math.log10((x+1)/70)))

# change the exploration rate over time.
def select_explore_rate(x):
    return max(min_explore_rate, min(1.0, 1.0 - math.log10((x+1)/70)))
```

### Result

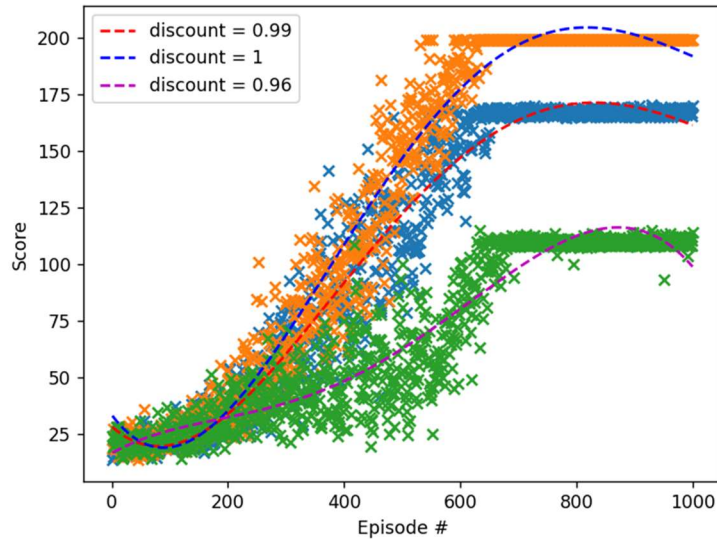


Figure 2 Plots of total reward received per episode (averaged over 10 trials) vs. number of episodes with three different discount rates.

Shown in Figure 2, as discount rate decreases, the average reward decreases. When discount = 1, 10 out of 10 trials succeed to solve the problem, whereas when discount = 0.96,

only 5 out of 10 trails solve the problem. Discount rates smaller than 0.96 hardly yield satisfactory results over 10 trails and therefore relative data are omitted here.

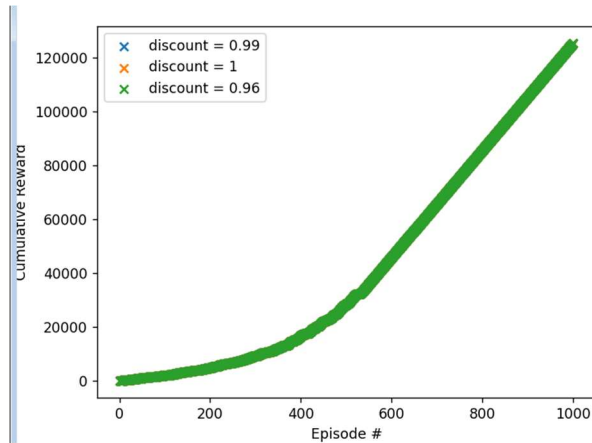


Figure 3 Plot of the cumulative reward (the sum of all rewards received so far) as a function of the number of episodes, with three different discount rates.

The overlapping results from seeding the random by 0 at the beginning of each trail. Seeding the pseudo-random number generators can help create reproductive result. Details are in task3Qlearning.py.

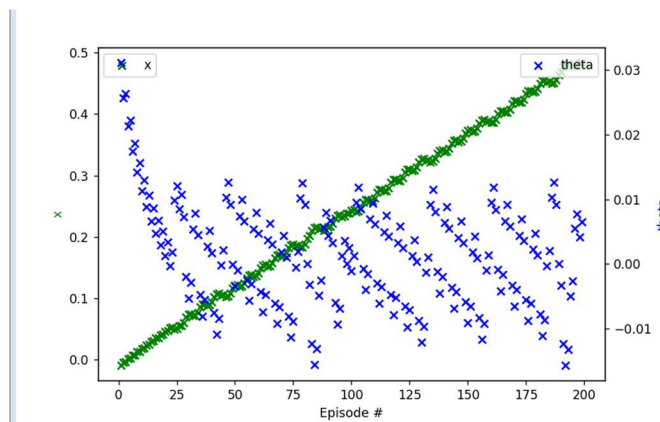


Figure 4 Plot of  $x$  and  $\theta$  as functions of time for episodes in a successful trail together with number of time steps during which the pole does not fail

Figure 4 shows  $x$  and  $\theta$  vs episodes in a successful trail. The cart is going towards one direction slowly but stay within the safe zone and the  $\theta$  is within a tiny range from 0.



## References

- [1] A. G. Barto, R. S. Sutton, and C. W. Anderson, "Neuronlike adaptive elements that can solve difficult learning control problems. SMC-13(5)," *IEEE Transactions on Systems, Man, and Cybernetics*, p. 834–846, 1983.
- [2] "Classic cart-pole system implemented by Rich Sutton et al.," [Online]. Available: [https://github.com/openai/gym/blob/master/gym/envs/classic\\_control/cartpole.py](https://github.com/openai/gym/blob/master/gym/envs/classic_control/cartpole.py).
- [3] "Practical Work: Reinforcement Learning," [Online]. Available: [https://scientific-python.readthedocs.io/en/latest/notebooks\\_rst/6\\_Machine\\_Learning/04\\_Exercices/02\\_Practical\\_Work/02\\_RL\\_1\\_CartPole.html](https://scientific-python.readthedocs.io/en/latest/notebooks_rst/6_Machine_Learning/04_Exercices/02_Practical_Work/02_RL_1_CartPole.html).
- [4] Barto, Sutton, and Anderson, "CartPole v0," [Online]. Available: <https://github.com/openai/gym/wiki/CartPole-v0>.