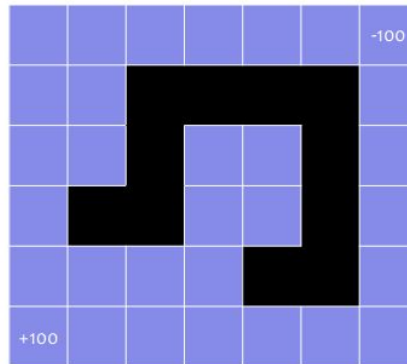


Reinforcement Learning

We extend the windy maze defined in P1 with probabilistic outcome after an action and two terminal states (left bottom and top right). It becomes a MDP problem. The maze map is shown in the following figure. However, we assume that the agent doesn't know either the reward function or the transition



model. The agent aims to run many trials in order to obtain Q-value for each (state, action) pair and the optimal action at each state.

Environment In your implementation, you need to simulate the windy maze environment: We assume that the wind comes from the south and the cost of one step for the agent is defined as follows: 1 for moving northward; 2 for moving westward or eastward; 3 for moving southward. The reward will be the negation of the reward. The agent can drift to the left or the right from the perspective of moving direction with probability 0.1. If the drifting direction is an obstacle, it will be bounced back to the original position. If the agent falls into any terminal state, it can't move out.

Reinforcement Learning In your implementation, you will generate many trials, each of which will result in a trajectory of (state, action, reward) tuple. The agent will use the ϵ -Greedy algorithm to choose an action at each state along each trajectory, where $\epsilon = 0.05$: the agent chooses a latest optimal action at each state with 95% and a random action with 5%. The initial state for each trial is chosen randomly and each trial will end at the goal state. Along each trajectory, the agent will use Q-Learning to update the Q-values. Since the reward function $R(s, a)$ here depends on both the state and the action taken at this state, the Q-value update equations should be revised accordingly.

$$N(s, a) \leftarrow N(s, a) + 1 \quad (1)$$

$$Q(s, a) \leftarrow Q(s, a) + \frac{1}{N_{s,a}} \left(R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a) \right) \quad (2)$$

We choose $\gamma = 0.9$.

Testing and Outputs In your testing, generate 20,000 trials starting from a random open square. We initialize the Q-values for each trial at any state-action as 0 except for two terminal states with +100 and -100 respectively. After 20,000 trials, report the following three outcomes for each algorithm:

- the access frequency at each state-action $N_{s,a}$;
- the Q-value function at each state-action $Q(s, a)$;
- the optimal action at each state-action.

The expected outcome should look like as follows:

• **Table of $N(s, a)$:**

1519		2062		3029		5021		8371		44	
48	50	98	90	3974	94	2382	100	1389	129	2013	29
1531		4474		70		101		148		29	
114		80								55	
85	99	5187	74	####		####		####		####	77 48
6595		75								4390	
115		112				5733		6770		3405	
129	123	2383	41	####		112 133		1083 97		####	697 1137
8654		38				1853		121		2056	
131						143		91		563	
202	185	####		####		74 54		1385 35		####	234 281
9680						4000		26		2612	
227		140		83		74				185	
221	221	9542	137	6461	81	4655	65	####		####	116 104
17401		122		92		71				4034	
		338		846		118		120		125	99
		7462	105	6616	88	6712	75	6103	77	5195	76 4340 63
		221		89		104		72		70	
										70	

• **Table of $Q(s, a)$:**

-5.2		-5.3		-5.7		-6.0		-6.5		-9.7	
3.8	-1.6	6.9	0.8	8.1	-3.6	1.6	-6.5	-2.9	-8.0	-7.3	-77.7 -100
24.7		17.0		-2.1		-6.0		-8.1		-11.7	
13.7		9.2								-69.1	
24.7	18.4	29.1	18.6	####		####		####		####	-14.5 -20.3
40.3		28.4								-12.2	
32.6		8.4				-6.0		-5.9		-10.7	
44.3	34.5	42.2	27.0	####		-3.0 -4.8		8.3 -6.0		####	-10.7 -10.7
54.5		30.8				17.4		-3.9		-10.1	
48.5						-0.3		-3.7		-9.2	
43.2	49.2	####		####		18.2 10.8		18.3 1.5		####	-9.1 -9.0
68.8						31.4		5.4		-5.4	
61.8		56.3		46.8		21.3				-7.0	
72.0	64.3	71.1	53.0	59.3	39.3	46.0	32.5	####		####	-5.5 -5.4
82.8		67.6		52.2		37.7				1.2	
		60.2		45.5		33.3		20.7		7.0	
+100		83.7	61.6	69.7	43.2	53.7	28.8	37.4	15.9	22.9	5.4 10.5 1.7
		55.7		55.5		38.9		25.9		11.0	
										4.2	

- Table of the optimal action at each state:

vvvv	vvvv	<<<<	<<<<	<<<<	<<<<	-100
vvvv	<<<<	####	####	####	####	vvvv
vvvv	<<<<	####	vvvv	<<<<	####	vvvv
vvvv	####	####	vvvv	<<<<	####	vvvv
vvvv	<<<<	<<<<	<<<<	####	####	vvvv
+100	<<<<	<<<<	<<<<	<<<<	<<<<	<<<<

where <<<<: moving westward; ^^^^: moving northward; >>>>: moving eastward; VVVV: moving southward; +100 and -100: the terminal rewards.

E-Greedy

Exploitation vs. Exploration

- The algorithm was split into 2 main components
 - Exploitation (95%) vs. Exploration (5%) selection
 - Desired direction (80%) vs. Drift probability each side (10% per side)
- From the initial random choice, we select either exploitation or exploration:

Exploitation Path

```
322 int eGreedy(int x, int y, int epsilonPercent=5)
323 {
324     //RNG check: 95% optimal action; 5% random action for epsilonPercent=5
325     int choice = rng(100);
326     int threshold = 100 - epsilonPercent;
327     int action = 0;
328
329     //Determine the action to be taken:
330     #pragma region Action determination
331     //From 1-95: Optimal action
332     if (choice <= threshold)
333     {
334         //Optimal action chosen
335         //Check for single highest Q-Val and execute the action
336         //If more than 1 highest Q-val, random between them
337         float wQVal = 0.0;
338         float nQVal = 0.0;
339         float eQVal = 0.0;
340         float sQVal = 0.0;
341         //Check west q-val: Col-1
342         wQVal = std::stof(qGrid[scaleGrid(x)][scaleGrid(y) - 1]);
343         //Check north q-val: Row-1
344         nQVal = std::stof(qGrid[scaleGrid(x) - 1][scaleGrid(y)]);
345         //Check east q-val: Col+1
346         eQVal = std::stof(qGrid[scaleGrid(x)][scaleGrid(y) + 1]);
347         //Check south q-val: Row+1
348         sQVal = std::stof(qGrid[scaleGrid(x) + 1][scaleGrid(y)]);
349
350         std::cout << "e-Greedy: {Exploit}\n";
351         //Action selected output: W/N/E/S
352         action = maxQSelect(wQVal, nQVal, eQVal, sQVal);
353     }
```

- In the event of selecting optimal routes, we select the action that has the highest Q-Value from the current grid position

- To select the proper action, and to account for potential equal Q-values that are also the highest, the maxQSelect function is used

```

225 int maxQSelect(float west, float north, float east, float south)
226 {
227     int action = 0;
228     bool wAction = false, nAction = false, eAction = false, sAction = false;
229     float maxQVal = 0;
230     int maxQCount = 0; //If there are equal max Q, keep track of equally optimal choices
231     maxQVal = std::max(west, std::max(north, std::max(east, south)));
232     float qVal[4] = { west, north, east, south };
233
234     for (int arrI = 0; arrI < 4; arrI++)
235     {
236         if (floatEquals(maxQVal, qVal[arrI]))
237         {
238             //Current direction is maxQVal
239             maxQCount++;
240             if (arrI == 0) wAction = true;
241             else if (arrI == 1) nAction = true;
242             else if (arrI == 2) eAction = true;
243             else if (arrI == 3) sAction = true;
244         }
245     }
246
247     //Action determination
248     if (maxQCount == 1)
249     {
250         //Solo maxQVal action selected
251         if (wAction == true) action = 1;
252         else if (nAction == true) action = 2;
253         else if (eAction == true) action = 3;
254         else if (sAction == true) action = 4;
255     }
256     else if (maxQCount == 4)
257     {
258         action = rng(4); //1 of 4 action selected
259     }

```

- For cases with only 1 max Q-value, the choice is straightforward
- Similarly, for cases with all Q-value being equal, a random action is selected

```

260     else
261     {
262         //from 2 or 3 specific actions to choose
263         //Create rng buckets for each maxQVal action
264         int bucket = rng(maxQCount * 100);
265         //std::cout << "Bucket rng: " << bucket << " | ";
266
267         //Each block of 100 numbers assigned to 1 action
268         if (maxQCount == 2)
269         {
270             if (bucket <= 100)
271             {
272                 if (wAction == true) action = 1;
273                 else if (wAction != true && nAction == true) action = 2;
274                 else if (wAction != true && nAction != true) action = 3;
275             }
276             else if (bucket > 100)
277             {
278                 if (wAction == true)
279                 {
280                     if (nAction == true) action = 2;
281                     else if (eAction == true) action = 3;
282                     else if (sAction == true) action = 4;
283                 }
284                 else if (nAction == true)
285                 {
286                     if (eAction == true) action = 3;
287                     else if (sAction == true) action = 4;
288                 }
289                 else
290                 {
291                     action = 4;
292                 }
293             }
294         }

```

- For 2 max Q-Val selection, we select the respective actions
- The bucket will generate numbers from 1-200 in maxQCount == 2 case
 - 1-100 would select the lowest possible action (ie: action 1-West, if that was available; if action 1-West was **not** the highest, it would check if action 2-North was available...)
 - 101-200 would select action 2 (similarly, the other possible action is the highest possible action available - action 4-South if it was indeed the highest Q-Val, and checks lower actions if not)
- For reference:
 - Action 1: West
 - Action 2: North
 - Action 3: East
 - Action 4: South

```

295     else if (maxQCount == 3)
296     {
297         if (bucket <= 100)
298         {
299             if (wAction == true) action = 1;
300             else if (wAction != true && nAction == true) action = 2;
301         }
302         else if (bucket > 100 && bucket <= 200)
303         {
304             if (wAction == true)
305             {
306                 if (nAction == true) action = 2;
307                 else if (nAction != true && eAction == true) action = 3;
308             }
309             else if (wAction != true) action = 3;
310         }
311         else
312         {
313             if (eAction == true) action = 3;
314             else if (sAction == true) action = 4;
315         }
316     }
317 }
318 //std::cout << "Action: " << action << "\n";
319 return action;
320 }

```

- Similarly, for 3 max Q-Val selection, the same process of selecting the proper respective actions as 2 max Q-Val case

Floating point precision

- Because floating points precision is an issue with "==" operator, a function with a cutoff resolution is used instead
- For practical purposes, differences smaller than the resolution ("epsilon") is considered equal

```

39 bool floatEquals(float a, float b)
40 {
41     //Function to prevent floating point comparison using "==" to be inaccurate due to "fuzziness"
42     float epsilon = 0.0000001;
43     //For practical purposes, differences smaller than epsilon would be considered "equal"
44     return std::abs(a - b) < epsilon;
45 }

```

Exploration Path

```
354 //From 96-100: Random action
355 else if (choice > threshold)
356 {
357     //Random action chosen (1 of 4 possible actions)
358     int randChoice = rng(4);
359     if (randChoice == 1) action = 1; //West
360     else if (randChoice == 2) action = 2; //North
361     else if (randChoice == 3) action = 3; //East
362     else if (randChoice == 4) action = 4; //South
363     std::cout << "e-Greedy: |Explore|\n";
364 }
365 std::cout << "Current position: " << x << ", " << y << "\n";
366 #pragma endregion Action determination
```

- A random action chosen with no regards to Q-Value
- From the action determination section, the “absolute” action is chosen that would be chosen in the event of 0 drift probability
- The next section will simply take the action into consideration and apply drift probability for the final “observed” action of the state-action pair

Drift Probability

```
368 #pragma region Drift
369 //Once maxQSelect or rng selects the desired direction, account for drift:
370 //80% to go towards desired action
371 //10% to drift to each side
372
373 int driftRng = rng(100);
374 //If 1-80: Desired direction | 81-90: Drift to sideA | 91-100: Drift to sideB
375 bool desiredDirection = false, driftA = false, driftB = false;
376 if (driftRng <= 80) desiredDirection = true;
377 else if (driftRng > 80 && driftRng <= 90) driftA = true;
378 else if (driftRng > 90) driftB = true;
379
380 std::cout << "Desired Action: " << action;
381 switch (action)
382 {
383 case 1://Desired direction is West
384     if (desiredDirection == true) action = 1; //West
385     else if (driftA == true) action = 2; //North
386     else if (driftB == true) action = 4; //South
387     break;
388 case 2://Desired direction is North
389     if (desiredDirection == true) action = 2; //North
390     else if (driftA == true) action = 1; //West
391     else if (driftB == true) action = 3; //East
392     break;
393 case 3://Desired direction is East
394     if (desiredDirection == true) action = 3; //East
395     else if (driftA == true) action = 2; //North
396     else if (driftB == true) action = 4; //South
397     break;
398 case 4://Desired direction is South
399     if (desiredDirection == true) action = 4; //South
400     else if (driftA == true) action = 1; //West
401     else if (driftB == true) action = 3; //East
402     break;
403 }
404 #pragma endregion Drift
405 if (desiredDirection != true) std::cout << " | Drift occurred | ";
406 std::cout << " | Actual Action: " << action << "\n";
407 //Final resulting action (takes into account maxQVal, rng, and drift)
408 //The true position that will end up happening
409 return action;
410 }
```

- Action determination from above is the “absolute” direction we would taken if there was no drift possibility
- Desired direction (action) is taken at a 80% rate, 10% for each side
- The returned action from e-Greedy is the actual observed action of the state-action pair (s,a)

Q-Learning

```
422 std::pair<int,int> qLearning(int x, int y, int epsilonPercent=5, int actionOverride=0)
423 {
424     int oldNVal = 0, newNVal = 0;
425     float oldQVal = 0.0, newQVal = 0.0;
426     float nextQVal = 0.0, deltaA = 0.0, deltaB = 0.0;
427     float discount = 0.9, actionReward = 0.0;
428     int action = 0;
429     if (actionOverride != 0)
430     {
431         action = actionOverride;
432     }
433     else
434     {
435         action = eGreedy(x, y, epsilonPercent);
436     }
437
438     int nextRow = 0, nextCol = 0;
439     float nextWQ = 0.0, nextNQ = 0.0, nextEQ = 0.0, nextSQ = 0.0;
440     int i = 0, j = 0;
441     //Input the state (position), eGreedy selects the action
442     switch (action)
443     {
444     case 1: //West (Col-1) |  $R(s,a) = -2$ 
445         i = scaleGrid(x), j = scaleGrid(y) - 1;
446         actionReward = -2.0;
447         if (grid[x][y - 1] == "####" || y - 1 < 0) { nextCol = y; }
448         else { nextCol = y - 1; }
449         nextRow = x;
450         break;
451     case 2: //North (Row-1) |  $R(s,a) = -1$ 
452         i = scaleGrid(x) - 1, j = scaleGrid(y);
453         actionReward = -1.0;
454         if (grid[x - 1][y] == "####" || x - 1 < 0) { nextRow = x; }
455         else { nextRow = x - 1; }
456         nextCol = y;
457         break;
458     case 3: //East (Col+1) |  $R(s,a) = -2$ 
459         i = scaleGrid(x), j = scaleGrid(y) + 1;
460         actionReward = -2.0;
461         if (grid[x][y + 1] == "####" || y + 1 >= col) { nextCol = y; }
462         else { nextCol = y + 1; }
463         nextRow = x;
464         break;
465     case 4: //South (Row+1) |  $R(s,a) = -3$ 
466         i = scaleGrid(x) + 1, j = scaleGrid(y);
467         actionReward = -3.0;
468         if (grid[x + 1][y] == "####" || x + 1 >= row) { nextRow = x; }
469         else { nextRow = x + 1; }
470         nextCol = y;
471         break;
472     }
```



```

474 //Update the state-action
475 std::pair<int, int> updatedPosition(nextRow, nextCol);
476
477 oldNVal = std::stoi(countGrid[i][j]);
478 newNVal = oldNVal + 1;
479
480 oldQVal = std::stof(qGrid[i][j]);
481
482 //Q-Learning update cases: Terminal or non-terminal states
483 if (endState(updatedPosition))
484 {
485     //Goal
486     if (grid[updatedPosition.first][updatedPosition.second] == "+100")
487     {
488         nextQVal = 100.0;
489     }
490     //Trap
491     else if (grid[updatedPosition.first][updatedPosition.second] == "-100")
492     {
493         nextQVal = -100.0;
494     }
495 }
496

```

- Keep track of $Q(s', a')$ in the event of terminal state
- This is more so for handling the way data is stored/retrieved from the larger scaled grid

```

497 else //Not terminal state
498 {
499     nextWQ = std::stof(qGrid[scaleGrid(nextRow)][scaleGrid(nextCol) - 1]);
500     nextNQ = std::stof(qGrid[scaleGrid(nextRow) - 1][scaleGrid(nextCol)]);
501     nextEQ = std::stof(qGrid[scaleGrid(nextRow)][scaleGrid(nextCol) + 1]);
502     nextSQ = std::stof(qGrid[scaleGrid(nextRow) + 1][scaleGrid(nextCol)]);
503
504     nextQVal = std::max(nextWQ, std::max(nextNQ, std::max(nextEQ, nextSQ)));
505 }
506 deltaA = 1.0 / (float)newNVal; // 1/N(s,a)
507 deltaB = actionReward + (discount * nextQVal) - oldQVal; //R(s,a) + (gamma * maxNextQVal) - oldQVal
508 newQVal = oldQVal + (deltaA*deltaB);
509
510 //Store iteration calculations
511 countGrid[i][j] = std::to_string(newNVal);
512 qGrid[i][j] = std::to_string(newQVal);
513
514 //std::cout << "Current (x,y): " << x << ", " << y << "\nAction: " << action << "\nNext (x,y): " << next
515 return updatedPosition;
516 }

```

- $N(s,a)$ has been stored using int, so (float) cast is necessary for proper delta calculations
- Once calculations are completed and the corresponding $N(s,a)$ and $Q(s,a)$ updated, return new position to begin the next action from the new state

e-Greedy: (Exploit)
 Current position: 4,1
 Desired Action: 1 | Actual Action: 1
 N(s,a):

Zishen Teh

24	182	118	77	66	29	
206	170	1124	15	840	14	632
1192	162	110	73	53	24	
27	79					3
229	264	534	8	####	####	####
1688	76					214
35	62		156	291		8
292	276	505	8	####	79	93
2148	62		442	44		416
42			22	60		12
378	283	####	####	113	104	297
2364			851	41		618
43	208	85	19			12
501	499	1468	31	660	12	152
3509	199	73	1211			985
	330	341	326	211	172	157
	2535	31	2496	37	2580	37
	338	378	371	222	168	138

9:09 PM
4/5/2020

Q(s,a):

39.33	35.08	29.67	24.8	20.42	16.48	
39.88	33.65	40.09	26.63	34.00	21.96	28.67
46.76	39.72	27.67	22.8	18.42	14.48	
39.53	35.08					-90.99
47.55	40.82	47.76	24.87	####	####	####
55.29	47.66					18.4
47.34	41.99		-2.93	-0.91		13.49
56.09	48.48	56.29	48.66	####	23.63	11.13
64.77	47.66			39.95	21.39	23.78
55.91				26.66	12.86	18.62
65.59	65.53	####	####	40.1	31.8	41.37
75.3				48.9	31.16	29.76
65.19	67.67	59	40.58			23.55
76.14	66.53	76.3	58	66.67	49.96	57.6
87	76.2	66.48	57.73			36.4
	67.25	58.83	58.64	51.86	44.77	31.34
+100	88	65.24	77.2	57.09	67.48	49.43
	76.2	66.48	57.73	49.86	42.77	36.4

Zishen Teh

9:11 PM
4/5/2020

```

Optimal action:
vvvv <<<< <<<< <<<< <<<< <<<< -100

vvvv <<<< ##### ##### ##### ##### vvvv
vvvv <<<< ##### vvvv <<<< ##### vvvv
vvvv ##### ##### vvvv <<<< ##### vvvv
vvvv <<<< <<<< vvvv ##### ##### vvvv
+100 <<<< <<<< <<<< <<<< <<<< <<<<

e-Greedy: {Exploit}
Current position: 4,0
Desired Action: 4 | Actual Action: 4

```

- Once the Q-Learning completes, the e-Greedy begins again from the new current position (4,0) - previously (4,1)
- E-greedy selects exploitation in this case, so the desired action is 4 (south)
- Because drift did not occur, the actual action taken is indeed south and reaches terminal state, which will end this specific trial run

```

e-Greedy: {Exploit}
Current position: 5,2
Desired Action: 1 | Drift occurred | Actual Action: 4
Q(s,a):

```

24	183	122	78	68	29						
207	172	1139	18	854	14	640	11	430	6	210	4
1206	165	110	76	53	24						
28	81				3						
233	266	540	8	####	####	####	####	29	33		
1710	78				220						
35	62	156	291	8							
293	279	514	8	####	79	96	325	45	####	63	67
2179	64	447	45	425							
43	22	60	12								
384	287	####	####	114	104	302	11	####	107	101	
2400		865	41	631							
45	218	89	19	12							
512	507	1490	32	677	12	150	185	####	####	136	113
3576	203	76	1230	1000							
	335	351	332	215	173	159					
2582	32	2538	38	2625	37	1522	26	1206	19	1099	18
346		389	376	228	170	141					

Zishen Teh

9:16 PM
4/5/2020

- From position (5,2), the desired action is 1 (West)
- However, drift occurred and the action ended up drifting to the south instead

