

FACULTAD DE CIENCIAS  
FÍSICAS Y MATEMÁTICAS  
UNIVERSIDAD DE CHILE

## EL7021-1 Seminario de Robótica y Sistemas Autónomos

---

### Entrega Tarea 2: Q-Learning y DQN

Fecha de entrega: Martes 26 de abril, 23:59 hrs.

Estudiante: Francisco Molina L.  
Profesor: Javier Ruiz del Solar  
Auxiliar: Francisco Leiva C.  
Semestre: Otoño 2022

## 1. Parte I: Q-Learning

1. En este problema el MDP es parcialmente observable, dado que el agente no tiene acceso directo a la información del entorno, por lo que no se conoce la matriz de transición de estados. Con esto en mente, el **POMDP**  $\langle \mathbf{S}, \mathbf{A}, \mathbf{R}, \gamma \rangle$  del carrito unidimensional se define por:

- El espacio de estados  $\mathbf{S}$ : Corresponde a valores continuos de la posición y velocidad del carrito, ambos con límites superiores e inferiores:

$$\mathbf{S}(x, v), \quad x \in [-1.2, 0.6], \quad v \in [-0.07, 0.07] \quad (1)$$

La posición  $x$  está delimitada por los límites espaciales de la simulación, en  $-1.2$  se sale del mapa por la izquierda y en  $0.6$  se llega al objetivo. Mientras que la velocidad  $v$  está delimitada por la potencia simulada del carrito, no puede ir más rápido que  $0.07$  en ambas direcciones.

- Las acciones  $\mathbf{A}$ : Corresponde a la aceleración que puede escoger el carrito, representada por números enteros del 0 al 2:
    - 0) Aceleración negativa, el carrito se mueve a la izquierda
    - 1) Aceleración neutra, el carrito se mueve por inercia
    - 2) Aceleración positiva, el carrito se mueve a la derecha
  - La función de recompensa  $\mathbf{R}$ : Corresponde a la recompensa que recibe el agente por ejecutar cada acción. Para que el agente aprenda a llegar a la meta de la manera más rápida posible se aplica como recompensa una penalización de  $-1.0$  por cada transición, excepto cuando se alcanza el objetivo donde se obtiene una recompensa de  $0.0$ .
  - El factor de descuento  $\gamma$ : Representa la incertidumbre sobre el futuro para el agente, hace que se valoren más o menos las recompensas futuras dependiendo del valor escogido y se escoge de manera arbitraria como  $0.95$ .
2. Para poder aplicar Q-Learning tabular al problema el espacio de estados debe ser discreto. Considerando que es un problema pequeño se puede discretizar la posición cada  $0.1$  unidades y la velocidad cada  $0.01$ .

Con esto se obtienen 19 valores posibles para la posición:

$$x \in [-1.2, -1.1, -1.0, -0.9, -0.8 \dots, 0.2, 0.3, 0.4, 0.5, 0.6] \quad (2)$$

Y 15 valores posibles para la velocidad:

$$v \in [-0.07, -0.06, -0.05, -0.04, \dots, 0.04, 0.05, 0.06, 0.07] \quad (3)$$

Por lo que el espacio de estados discretizado  $\mathbf{S}[x, v]$  corresponde a las 285 combinaciones posibles entre los  $x$  y  $v$  recién descritos.

3. La función `select_action` recibe una observación de la cual se extrae el estado  $s$  y devuelve:

$$\text{select\_action} \rightarrow \begin{cases} \text{La acción } a \text{ que maximiza } Q(s, a) & \text{con probabilidad } \epsilon \\ \text{Una acción } a \text{ aleatoria} & \text{con probabilidad } 1 - \epsilon \end{cases}$$

La función `update` recibe una tupla  $(s_t, s_{t+1}, a, r)$  y si  $s_t$  no es un estado terminal, actualiza su  $Q$ -Value según la regla:

$$Q^\pi(s_t, a_t) \leftarrow Q^\pi(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q^\pi(s_{t+1}, a') - Q^\pi(s_t, a_t)) \quad (4)$$

4. Usando los parámetros por defecto se obtiene la recompensa promedio y la tasa de éxito presentadas en las **Figuras 1 y 2**:

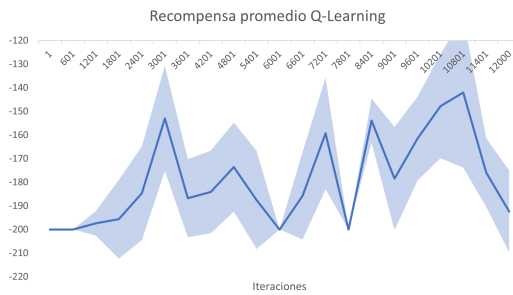


Figura 1: Recompensa promedio



Figura 2: Tasa de éxito

5. Implementando  $\epsilon$ -greedy se obtiene la recompensa promedio y la tasa de éxito presentadas en las **Figuras 3 y 4**:

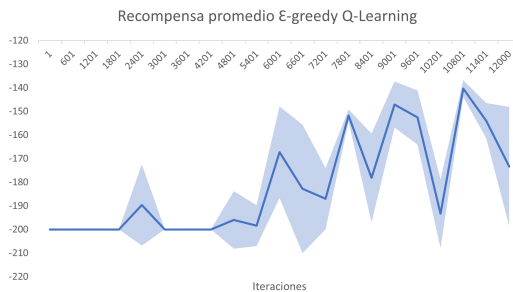


Figura 3: Recompensa promedio con  $\epsilon$ -greedy

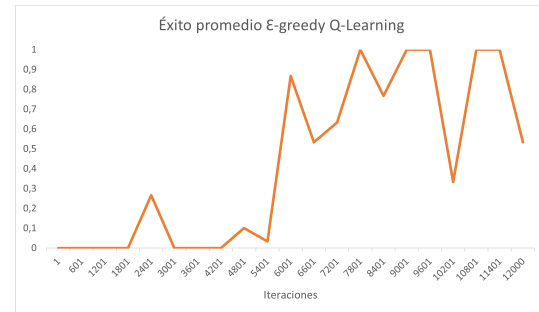


Figura 4: Tasa de éxito con  $\epsilon$ -greedy

Se puede apreciar que al implementar  $\epsilon$ -greedy se obtienen resultados mejores y más consistentes, especialmente en las últimas iteraciones. En la **Figura 4** se ve claramente que cerca de las 5400 iteraciones la tasa de éxito se dispara ( $\epsilon$  está llegando a 0) y se mantiene alrededor de 0.8 de ahí en adelante.

## 2. Parte II: Deep Q-Network

### 1. Trabajando en el archivo `deep_qnetwork.py`:

1.1. Para completar la clase `DeepQNetwork` basta con definir una red neuronal con dos capas ocultas de 64 neuronas, usando funciones de activación ReLU:

```
1 class DeepQNetwork(nn.Module):
2
3 def __init__(self, dim_states, dim_actions):
4     super(DeepQNetwork, self).__init__()
5     self.device = torch.device("cuda" if torch.cuda.is_available()
6                               else "cpu")
7     self.fc1 = nn.Linear(dim_states, 64)
8     self.fc2 = nn.Linear(64, 64)
9     self.fc3 = nn.Linear(64, dim_actions)
10    self.float()
11
12 def forward(self, input):
13     out = input.to(self.device)
14     out = F.relu(self.fc1(out))
15     out = F.relu(self.fc2(out))
16     out = self.fc3(out)
17     return out
18
```

Listing 1: Clase `DeepQNetwork`

1.2. Para completar el constructor de la clase `DeepQNetworkAgent` se crea la *target network* copiando la red `deep_QNetwork` y se inicializan el optimizador y la pérdida

```
1 # Complete
2 self._deep_QNetwork = DeepQNetwork(self._dim_states,
3                                    self._dim_actions)
4 self._target_deepQ_network = copy.deepcopy(self._deep_QNetwork)
5
6 # Adam optimizer and MSE Loss
7 self._optimizer = torch.optim.Adam(self._deep_QNetwork.parameters(),
8                                     lr=lr)
9 self._criterion = nn.MSELoss()
10
```

Listing 2: Clase `DeepQNetworkAgent`

- 1.3. Para completar el método `select_action` se devuelve la acción que maximiza la salida de la red `deep_QNetwork` ( $\text{argmax}_{a'} Q_{\theta}(s, a')$ ) con probabilidad  $\epsilon$  o una acción aleatoria con probabilidad  $1 - \epsilon$ . Además, se reduce el valor de  $\epsilon$  con cada iteración:

```
1  def select_action(self, observation, greedy=False):
2      if (np.random.random() > self._epsilon) or greedy:
3          with torch.no_grad():
4              Q_values = self._deep_QNetwork(torch.tensor(observation))
5              action = Q_values.argmax()
6      else:
7          action = torch.tensor(np.random.randint(0, self._dim_actions))
8
9      if not greedy and self._epsilon >= self._epsilon_min:
10         self._epsilon -= self._epsilon_decay
11     return action
12
```

Listing 3: `select_action`

## 2. Trabajando en el archivo `replay_buffer.py`:

2.1. Para completar el constructor de la clase `ReplayBuffer` se inicializan los *buffer* de  $s_t$ ,  $a_t$ ,  $r_t$ ,  $s_{t+1}$  y *done* como *arrays* llenos de ceros con sus respectivas formas:

```

1  def __init__(self, dim_states, dim_actions, max_size, sample_size):
2      assert sample_size < max_size, "Sample size cannot be greater than
3          buffer size"
4
5      self._buffer_idx = 0
6      self._exps_stored = 0
7      self._buffer_size = max_size
8      self._sample_size = sample_size
9
10     self._s_t_array = np.zeros((self._buffer_size, dim_states))
11     self._a_t_array = np.zeros((self._buffer_size, 1))
12     self._r_t_array = np.zeros((self._buffer_size, 1))
13     self._s_t1_array = np.zeros((self._buffer_size, dim_states))
14     self._term_t_array = np.zeros((self._buffer_size, 1))
15

```

Listing 4: Clase `ReplayBuffer`

2.2. Para completar el método `store_transition` simplemente se añaden los elementos a los buffer según `_buffer_idx`, se aumenta en 1 la cantidad de elementos almacenados y `_buffer_idx` se aumenta de modo que se reinicie cuando alcance el tamaño máximo del buffer:

```

1  def store_transition(self, s_t, a_t, r_t, s_t1, done_t):
2      # Add transition to replay buffers
3      self._s_t_array[self._buffer_idx] = s_t
4      self._a_t_array[self._buffer_idx] = a_t
5      self._r_t_array[self._buffer_idx] = r_t
6      self._s_t1_array[self._buffer_idx] = s_t1
7      self._term_t_array[self._buffer_idx] = done_t
8
9      # Update replay buffer index
10     self._buffer_idx = (self._buffer_idx + 1) % self._buffer_size
11     self._exps_stored += 1
12

```

Listing 5: `store_transition`

- 2.3. Para completar el método `sample_transitions` se genera una lista de índices aleatorios, con valores desde 0 hasta `_buffer_idx` si el *buffer* no está lleno, o con valores hasta `_buffer_size` en caso contrario. Luego simplemente se devuelven los elementos de cada *buffer* correspondientes a los índices aleatorios:

```

1  def sample_transitions(self):
2      assert self._exps_stored + 1 > self._sample_size
3      if self._exps_stored < self._sample_size:
4          sample_ids = random.sample(range(0, self._buffer_idx),
5                                     self._sample_size)
6
7      else:
8          sample_ids = random.sample(range(0, self._buffer_size),
9                                     self._sample_size)
10
11     return (self._s_t_array[sample_ids],
12            self._a_t_array[sample_ids].astype(int),
13            self._r_t_array[sample_ids],
14            self._s_t1_array[sample_ids],
15            self._term_t_array[sample_ids])
16

```

Listing 6: `sample_transition`

- 2.4. Para probar la efectividad del `ReplayBuffer` se usa el mismo ambiente ‘`CartPole-v1`’ con un *buffer* de tamaño 20 y muestras de tamaño 10. Usando la herramienta *debugger* de cualquier IDE, se pueden analizar las variables cuando se ejecuta el código. Con esto se puede apreciar que el `ReplayBuffer` efectivamente almacena la primera transición, en *arrays* de largo 20, como se ve en al **Figura 5**:

```

> self._a_t_array = (ndarray: (20, 1)) [[1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.] ...View as Array
> self._r_t_array = (ndarray: (20, 1)) [[1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.] ...View as Array
> self._s_t1_array = (ndarray: (20, 4)) [[ 0.03114567  0.17200139  0.00621174 -0.24465623], [ 0.  0.  0.  0. ], [ 0.  0.
> self._s_t_array = (ndarray: (20, 4)) [[ 0.03160656 -0.02304438  0.00528464  0.04635467], [ 0.  0.  0.  0. ], [ 0.  0.
> self._term_t_array = (ndarray: (20, 1)) [[0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.] ...View as Array

```

Figura 5: `ReplayBuffer` almacenando primera transición

Luego, cuando se ejecuta la línea:

```

1  st, a, r, st1, end = self.replay_buffer.sample_transitions()

```

Se puede ver que efectivamente el `ReplayBuffer` devuelve muestras como *arrays* de largo 10, como se ve en la **Figura 6**:

```

> st_arr = (ndarray: (10, 4)) [[ 0.09779516  0.56583208 -0.09386636 -0.91243911], [ 0.  0.  0.  0. ], [ 0.  0.  0.  0. ],
> a_arr = (ndarray: (10, 1)) [[0., 0., 0., 0., 0., 1., 1., 0., 0., 0.] ...View as Array
> r_arr = (ndarray: (10, 1)) [[1., 0., 0., 0., 0., 1., 1., 1., 1., 0.] ...View as Array
> st1_arr = (ndarray: (10, 4)) [[ 0.1091118  0.37209684 -0.11211514 -0.65067172], [ 0.  0.  0.  0. ], [ 0.  0.  0.  0. ],
> end_arr = (ndarray: (10, 1)) [[0., 0., 0., 0., 0., 0., 0., 0., 0., 0.] ...View as Array

```

Figura 6: `ReplayBuffer` muestreando décima transición

En el caso en el que el *buffer* está vacío o tiene menos elementos que *\_sample\_size* se puede ver que sólo se pueden almacenar transiciones pero no muestrear, pues para obtener la **Figura 6** se tuvieron que correr 10 iteraciones primero. Por otra parte, si *\_sample\_size* es mayor o igual al tamaño máximo del *buffer*, el script no corre gracias al *assert* definido en *sample.transitions*.

### 3. Continuando con *deep\_qnetwork.py*:

- 3.1. Para completar el método *replace\_transition\_network* basta con simplemente cargar los pesos de la red *deep\_QNetwork* hacia la *target network*:

```
1 def replace_target_network(self):
2     self._target_deepQ_network.load_state_dict(self._deep_QNetwork.
3                                                 state_dict())
4
```

Listing 7: *replace\_target\_network*

- 3.2. Por último, para completar el método *update* se sigue una secuencia de pasos. Primero se muestrean los *batches* de datos y se transforman a tensores con sus respectivos *dtypes*:

```
1 def update(self):
2     # get batches
3     st_arr, a_arr, r_arr, stl_arr, end_arr = self.replay_buffer.
4                                             sample_transitions()
5
6     # make tensors from batches
7     st_arr = torch.from_numpy(st_arr).float()
8     a_arr = torch.from_numpy(a_arr).type(torch.int64)
9     r_arr = torch.from_numpy(r_arr).float()
10    stl_arr = torch.from_numpy(stl_arr).float()
11    end_arr = torch.from_numpy(end_arr).bool()
12
```

Luego, se calculan las salidas esperadas con el método auxiliar *compute\_y*, el cual computa los *Q-values* máximos según la *target network* y luego llena el tensor *y<sub>j</sub>* con *r<sub>j</sub>* si *end<sub>j</sub>* es True (estado terminal) y con *r<sub>j</sub>* más los *Q-values* calculados anteriormente si *end<sub>j</sub>* es False:

```
1 y_t = self.compute_y(r_arr, stl_arr, end_arr)
2     |
3     v
4 def compute_y(self, r, stl, end):
5     with torch.no_grad():
6         QValues = torch.amax(self._target_deepQ_network(stl), dim=1)
7         y_j = torch.where(end.flatten(), r.flatten(),
8                           torch.add(end.flatten(), QValues,
9                                     alpha=self._gamma))
10    return y_j
11
```



Se llevan los gradientes a cero y se calculan los  $Q$ -values de la red `deep_QNetwork` con el método auxiliar `compute_QValues`. El método simplemente propaga los estados  $s_t$  por la red `deep_QNetwork` y por cada salida se escoge el  $Q$ -value de la acción  $a_t$  perteneciente al mismo episodio que  $s_t$ :

```
1 self._optimizer.zero_grad(set_to_none=True)
2 Q_values = self.compute_QValues(st_arr, a_arr)
3         |
4         v
5 def compute_QValues(self, st, a):
6     QValues = self._deep_QNetwork(st)
7     return torch.gather(QValues, dim=1, index=a).flatten()
8
```

Por último, se calcula la *loss*, se propaga con *backpropagation*, se limitan los gradientes a  $(-1, 1)$  y se realiza un paso de optimización:

```
1 loss = self._criterion(Q_values, y_t)
2 loss.backward()
3 torch.nn.utils.clip_grad_norm_(self._deep_QNetwork.parameters(), 1)
4 self._optimizer.step()
5
```

4. Se fijan los parámetros en `nb_training_steps:20000`, `gamma:0.99`, `epsilon:0.8` y `lr:0.01`. Cabe destacar que en estos experimentos de manera adicional se almacenan los pesos de la red según si superan la máxima recompensa promedio obtenida y dichos pesos se cargan en la última evaluación del agente.

- 4.1. A continuación se presentan los resultados del **Experimento 1**, donde se analiza el efecto de variar el tamaño máximo del *buffer*:

exp\_11 (1000, 16, 100)

Al ejecutar el exp\_11 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 7 y 8**:

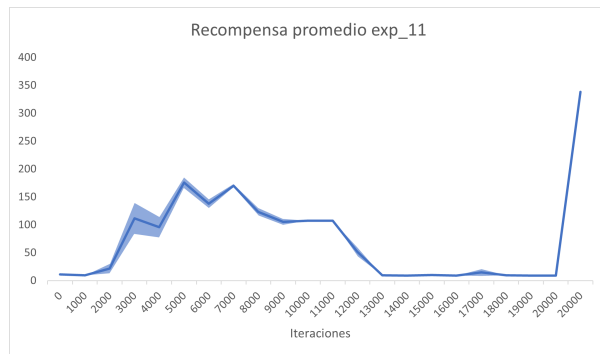


Figura 7: Recompensa promedio exp\_11



Figura 8: Tasa de éxito exp\_11

exp\_12 (2500, 16, 100)

Al ejecutar el exp\_12 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 9 y 10**:



Figura 9: Recompensa promedio exp\_12



Figura 10: Tasa de éxito exp\_12

exp\_13 (5000, 16, 100)

Al ejecutar el exp\_13 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 11 y 12**:



Figura 11: Recompensa promedio exp\_13



Figura 12: Tasa de éxito exp\_13

A continuación se presentan los resultados del **Experimento 2**, donde se analiza el efecto de variar el tamaño de las muestras:

exp\_21 (5000, 32, 100)

Al ejecutar el exp\_21 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 13 y 14**:



Figura 13: Recompensa promedio exp\_21



Figura 14: Tasa de éxito exp\_21

exp\_22 (5000, 64, 100)

Al ejecutar el exp\_22 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 15 y 16**:



Figura 15: Recompensa promedio exp\_22



Figura 16: Tasa de éxito exp\_22

exp\_23 (5000, 128, 100)

Al ejecutar el exp\_23 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 17 y 18**:



Figura 17: Recompensa promedio exp\_23



Figura 18: Tasa de éxito exp\_23

A continuación se presentan los resultados del **Experimento 3**, donde se analiza el efecto de variar la cantidad de iteraciones que se espera para actualizar la *target network*:

exp\_31 (5000, 128, 1)

Al ejecutar el exp\_31 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 19 y 20**:

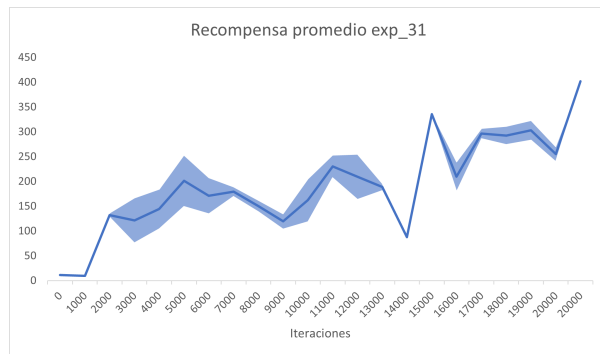


Figura 19: Recompensa promedio exp\_31



Figura 20: Tasa de éxito exp\_31

exp\_32 (5000, 128, 100)

Dados los parámetros escogidos el exp\_32 es el mismo experimento que el exp\_23, por lo que sus resultados y por ende su análisis son iguales. Es decir, para este experimento se deben considerar la recompensa promedio y la tasa de éxito presentados en las **Figuras 17 y 18**.

exp\_33 (5000, 128, 1000)

Al ejecutar el exp\_33 cinco veces y promediando los resultados se obtienen la recompensa promedio y la tasa de éxito promedio presentados en las **Figuras 21 y 22**:

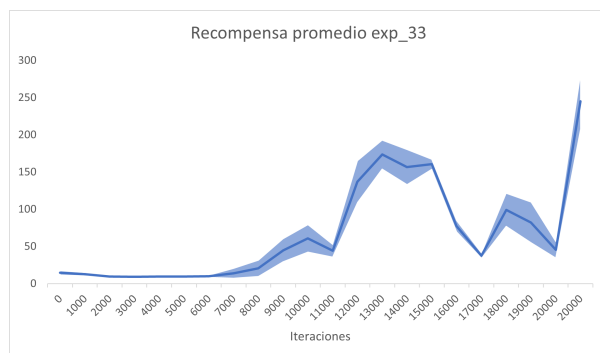


Figura 21: Recompensa promedio exp\_33



Figura 22: Tasa de éxito exp\_33

#### 4.2. A continuación se presenta el análisis de los experimentos recién mostrados.

El **Experimento 1** varía el tamaño máximo del *buffer* desde 1000 hasta 5000. Una idea central detrás de usar un *replay buffer* es evitar la correlación entre las muestras que le llegan a la red, para así mejorar su desempeño. Por ende, *buffers* más grandes promueven las chances de escoger muestras no correlacionadas, por lo que el aprendizaje de la red debiese ser más estable.

Si bien los resultados presentados en las **Figuras 7, 9 y 11** no logran los 500 puntos de recompensa, sí se condicen con lo recién descrito en el párrafo anterior. En la **Figura 7** el agente aprende un poco pero tiene una caída a las 12 mil iteraciones de la cual no se recupera. En la **Figura 9** el agente tiene una pérdida catastrófica también cerca de las 13 mil iteraciones. En cambio, en la **Figura 11** el agente obtiene recompensas, si bien no muy grandes, más estables. Estas bajas recompensas probablemente se deben a que el tamaño de las muestras es muy pequeño.

El **Experimento 2** varía el tamaño de las muestras desde 32 hasta 128. El tamaño de las muestras va de la mano con la idea del *replay buffer*. Muestras más grandes permiten que el agente se entrene con más episodios en menos tiempo, lo cual debiese acelerar su rendimiento. Sin embargo, tomar muchas muestras (casi tantas como el máximo) aumenta las chances de que estas por azar se encuentren correlacionadas.

Los resultados presentados en las **Figuras 13, 15 y 17** también se condicen con lo discutido anteriormente. Gracias al tamaño de su *buffer* no tienen muchas pérdidas catastróficas y se puede apreciar que a medida que aumenta el tamaño de las muestras los resultados también mejoran, obteniendo los mejores en el exp\_23 (**Figura 17**).

Por último, el **Experimento 3** varía la cantidad de iteraciones que se espera para actualizar los pesos de la *target network*. El propósito de incluir la *target network* es estabilizar el entrenamiento del agente. Si se actualiza la *target network* muy seguido, el objetivo hacia el cual se entrena la red se mueve mucho y por ende el entrenamiento es inestable. En el otro extremo, si la *target network* se actualiza muy poco la red nunca va a aprender y se quedará estancada con malos resultados. Por ende, hay que encontrar un “punto dulce” para la actualización de la *target network*, no muy bajo pero tampoco no muy grande.

Los resultados presentados en las **Figuras 19, 17 y 21** también se condicen con lo explicado en el párrafo anterior. En la **Figura 19** la red se actualiza en cada iteración, sorprendentemente obtiene buenos resultados, pero se puede apreciar que tiene bastantes cambios bruscos y que hay mucha varianza de un experimento a otro (mostrado por las sombras de error). En la **Figura 17** la red se actualiza cada 100 iteraciones, los resultados no tienen mucha varianza y aumentan consistentemente. En la **Figura 21** el agente se estanca y no logra aprender, llegando a una máxima tasa de éxito de 0.25 (**Figura 22**). Por ende, el “punto dulce” se encuentra con reemplazar la *target network* cada 100 iteraciones, siendo ese el mejor modelo para este problema.