

Construcción de Carro Autónomo mediante algoritmo DQN con Pytorch y Kivy.

Colorado, Cesar
Universidad Nacional
de Ingeniería
Lima, Perú
ccoloradoc@uni.pe

Inocente, Patricks
Universidad Nacional
de Ingeniería
Lima, Perú
pinocentev@uni.pe

Guevara, Junior
Universidad Nacional
de Ingeniería
Lima, Perú
jguevarat@uni.pe

Sulca, Leonardo
Universidad Nacional
de Ingeniería
Lima, Perú
leonardo.sulca.d@uni.pe

Resumen—Deep Q-network (DQN) es un algoritmo propio del aprendizaje por refuerzo profundo, que tiene como herramienta principal una red neuronal usada para la estimación del valor Q (función Valor Acción) de las acciones disponibles para un estado en particular. En este trabajo se explora el uso de DQN para resolver la problemática de un carro autónomo que debe viajar desde el centro de una ciudad hacia un aeropuerto y viceversa; por lo que en cada momento se necesita obtener una ruta óptima para su recorrido, esta experiencia es simulada usando la plataforma Kivy y la librería pytorch de python; de tal forma que Kivy actúe como visualizador y pytorch nos proporcione las herramientas para la implementación de nuestra DQN, dando como resultado varios agentes que logran viajar desde el origen hacia la meta evitando choques entre sí y sorteando obstáculos a medida que estos son instanciados en el mapa de forma interactiva.

Índice de Términos— DQN, Kivy, pytorch, Carro autónomo.

I. INTRODUCCIÓN

Kivy es principalmente un framework de Python gratuito y de código abierto para desarrollar aplicaciones multitáctiles con una interfaz de usuario natural. El objetivo principal de usar Kivy en este proyecto es obtener un entorno de visualización y de interacción dinámica para la puesta de obstáculos en el mapa. El mapa se compone de 3 zonas principales establecidas para el recorrido del carro autónomo (véase Figura 1).

1. **Centro de la Ciudad:** Zona considerada como punto de origen o meta, dependiendo del viaje actual que debe realizar nuestro carro autónomo.
2. **Aeropuerto:** Zona considerada también como origen o meta.
3. **Ciudad:** Describe la zona intermedia del mapa donde encontraremos: - **Obstáculos:** Líneas dibujadas en tiempo de ejecución que no deberían ser transitadas por el carro autónomo. - **Avenidas:** Zona libre de obstáculos.

II. ESTADO DEL ARTE

Q-learning fue introducido por *Watkins* en 1989. *Watkins* presentó una prueba de convergencia a *Dayan* en 1992.

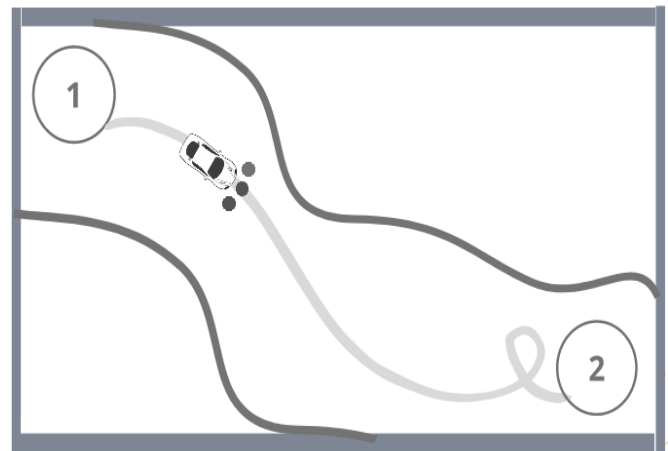


Figura 1. Representación del Entorno (Creación propia). Se muestra las partes del Mapa. 1. Centro de Ciudad. 2. Aeropuerto. La línea de color más suave es una ruta representativa que tomaría nuestro carro autónomo dado los obstáculos marcados con líneas de color más fuerte.

Una prueba matemática más detallada fue expuesta por *Tsitsiklis* en 1994, y por *Bertsekas* y *Tsitsiklis* en su libro de 1996 **Neuro-Dynamic Programming**.

Watkins lo señalaba ya el título de la tesis de su doctorado "Learning from delayed rewards" (Aprendiendo de las recompensas tardías). Ocho años antes en 1981 el mismo problema, bajo el nombre de "aprendizaje por refuerzo tardío", que fue solucionado por *Bozinovski's Crossbar Adaptive Array* (CAA).

La matriz de memoria $W(a,s)$ es exactamente igual que la Q-table de Q-learning de ocho años más tarde. La arquitectura introdujo el término "evaluación de estado" en aprendizaje por refuerzo. El algoritmo de aprendizaje crossbar (escrito en pseudocódigo matemático en papel) realiza los siguientes cálculos en cada iteración:

- En el estado s efectúa la acción a
- En consecuencia recibe el estado s'

- Realiza la valuación del estado $v(s')$
- Actualiza el valor crossbar $W'(a,s) = W(a,s) + v(s')$

El término “refuerzo secundario” viene de la teoría de aprendizaje animal, utilizada para modelar los valores de estado a través de propagación hacia atrás: el valor de estado $v(s')$ de la situación consecuente es propagado hacia atrás hacia las situaciones encontradas previamente. CAA calcula los valores de estados verticalmente y las acciones horizontalmente. Grafos de demostración nos muestran los estados contenidos en el aprendizaje por refuerzo tardío (deseables, indeseables, y estados neutros), que han sido calculados por la función de evaluación de estados. Este sistema de aprendizaje fue el precursor del algoritmo Q-learning.

En 2014 Google DeepMind patentó una aplicación de Q-learning de cara al aprendizaje profundo, llamado “aprendizaje de refuerzo profundo.” “Q-learning profundo” que puede jugar a los juegos Atari 2600 en niveles humanos expertos

III. VARIANTES:

Q-learning profundo:

El sistema DeepMind utilizó redes neuronales convolucionales profundas, con capas de filtros convolucionados enladrillados para asemejar los efectos de los campos receptivos. El aprendizaje por refuerzo es inestable o divergente cuando un aproximador de funciones no lineales se utiliza para representar Q. Esta inestabilidad proviene de la correlación presente en las secuencias de observación, el hecho de que actualizaciones pequeñas de Q puedan cambiar significativamente la política y la distribución de los datos, y las correlaciones entre Q y los valores objetivo.

La técnica utilizó repetición de experiencia, un mecanismo biológicamente inspirado que usa una muestra aleatoria de acciones previas en vez de la acción más reciente para continuar. Esto elimina las correlaciones en la secuencia de observación y suaviza los cambios en la distribución de los datos. La actualización iterativa ajusta Q de cara a los valores objetivo que solo se actualizan periódicamente, reduciendo más adelante las correlaciones con el objetivo.

Q-learning doble:

Puesto que el valor futuro máximo aproximado de una acción en Q-learning se evalúa utilizando la misma función Q que en la política de acción actualmente seleccionada, en entornos ruidosos Q-learning a veces puede sobrestimar los valores de la acción, retrasando el aprendizaje. Se propuso una variante llamada Q-learning doble para corregir esto. Q-learning doble es un algoritmo de aprendizaje por refuerzo sin política,

donde una política diferente se utiliza para la evaluación del valor que se usa para seleccionar próxima acción.

En la práctica, dos funciones de valor diferentes se entrenan en Q^A una moda simétrica mutua usando experiencias separadas, Q^A y Q^B . Q^A Q^B El paso de actualización de Q-learning doble es, por tanto, así:

$$Q_{t+1}^A(s_t, a_t) = Q_t^A(s_t, a_t) + \alpha_t(s_t, a_t) \left(r_t + \gamma Q_t^B(s_{t+1}, \arg \max_a Q_t^A(s_{t+1}, a)) - Q_t^A(s_t, a_t) \right)$$

, y Ahora el valor estimado del futuro descontado se evalúa utilizando una política diferente, lo cual soluciona el problema de la sobrestimación.

Este algoritmo se combinó más tarde con aprendizaje profundo, como el algoritmo DQN, resultando en doble DQN, el cual supera en rendimiento el algoritmo DQN original.

IV. ARTICULOS:

- Hester, T., Vecerik, M., Pietquin, O., Lanctot, M., Schaul, T., Piot, B., ... Gruslys, A. (2018, April). Deep q-learning from demonstrations. In Thirty-second AAAI conference on artificial intelligence.
- Van Hasselt, H., Guez, A., Silver, D. (2016, March). Deep reinforcement learning with double q-learning. In Proceedings of the AAAI conference on artificial intelligence (Vol. 30, No. 1).
- Ong, H. Y., Chavez, K., Hong, A. (2015). Distributed deep Q-learning. arXiv preprint arXiv:1508.04186.

V. DISEÑO DE EXPERIMENTO

Viendo los casos en que algunos equipos de nuestro grupo hacían los movimientos de cada carrito en diferente FPS, decidimos establecer un FPS estable para que detecte el obstáculo y aprenda correctamente.

```
1 Config.set('input', 'mouse', 'mouse,
    multitouch_on_demand')
2 Config.set('graphics', 'maxfps', '20')
```

Aquí haremos la inicialización del Modelo con los parámetros “input_size”, “nb_action”, “gamma”

```
1 class Dqn(object):
2
3     def __init__(self, input_size, nb_action, gamma):
4         :
5         self.start = time.time()
6         self.gamma = gamma
7         self.model = Network(input_size, nb_action)
8         self.memory = ReplayMemory(capacity=100000)
9         self.optimizer = optim.Adam(params=self.
    model.parameters())
10        self.last_state = torch.Tensor(input_size).
    unsqueeze(
11            0)
12        self.last_action = 0
```

```

12         self.last_reward = 0
13         self.reward_window = []

```

Aquí definimos el uso de una función para la toma de acciones, donde tendremos en cuenta la probabilidad de cada estado actual

```

1 def select_action(self, state):
2     probs = F.softmax(self.model(Variable(state)
3         )*100)
4     action = probs.multinomial(len(probs))
5     return action.data[0, 0]

```

Aquí definimos el uso de una función el aprendizaje:

```

1 def learn(self, batch_states, batch_actions,
2     batch_rewards, batch_next_states):
3     # Conseguimos el output del Modelo.
4     # Gather hace la eleccion, con unsqueeze
5     # agregamos una dimension para seguir
6     # trabajando con la eleccion, y despues con
7     # squeeze le bajamos a la dimension inicial
8     # para obtener las salidas del Modelo
9     batch_outputs = self.model(batch_states).
10    gather(
11        1, batch_actions.unsqueeze(1)).squeeze
12    (1)
13
14    # Recibe una tupla de dimension 3 y devuelve
15    # el maximo valor
16    # de cada una de las tuplas evaluadas, y se
17    # almacena en batch_next_outputs
18    # print(batch_next_outputs)
19    batch_next_outputs = self.model(
20    batch_next_states).detach().max(1)[0]
21
22    # Q(s,a) = R(s,a) + gamma*max(Q(s',a))
23    # Q(varios) para el lote examinado
24    batch_targets = batch_rewards + self.gamma *
25    batch_next_outputs
26
27    # Calculamos la perdida y tratamos de
28    # minimizar la perdida con el Optimizador
29    # Adam
30    td_loss = F.smooth_l1_loss(batch_outputs,
31    batch_targets)
32    self.optimizer.zero_grad() # -> Poner los
33    # gradientes a Cero
34    td_loss.backward() # -> Backpropagation,
35    # obtenemos las gradientes
36    self.optimizer.step() # -> actualizando los
37    # pesos con Adam

```

Aquí definimos el uso de una función para la actualización de un estado:

```

1 def update(self, new_state, new_reward):
2     # Convierte el estado a un tensor float, y
3     # le aumenta 1 dimension
4     new_state = torch.Tensor(new_state).float().
5     unsqueeze(0)
6
7     # Agregando a la memoria un evento, el
8     # evento es:
9     # El estado actual, la accion, la recompensa
10    # , y el estado siguiente

```

```

7     # last_state es el estado actual
8     # last_reward -> la recompensa inmediata
9     # new_state -> el estado siguiente
10    self.memory.push((self.last_state, torch.
11    LongTensor(
12        [int(self.last_action)]), torch.Tensor([
13        self.last_reward]), new_state))
14
15    # Le pasamos un estado, y con select_action
16    # elije la accion
17    # a tomar, de las probabilidades con la
18    # multinomial
19    # tomando una accion [0,1,2]
20    new_action = self.select_action(new_state)
21
22    # Si la memoria de eventos se llena
23    # significa que
24    # Empezamos a aprender tomando 100 muestras
25    # Y lo mandamos a aprender, siempre y cuando
26    # tengamos
27    # + de 100 muestras
28
29    if len(self.memory.memory) > 100:
30        batch_states, batch_actions,
31        batch_rewards, batch_next_states = self.memory.
32        sample(
33            100)
34        self.learn(batch_states, batch_actions,
35            batch_rewards,
36            batch_next_states)
37        self.last_state = new_state
38        self.last_action = new_action
39        self.last_reward = new_reward
40        self.reward_window.append(new_reward)
41        if len(self.reward_window) > 1000:
42            del self.reward_window[0]
43        return new_action

```

De nuestro archivo mapa.py

Función de actualización:

```

1 def update(self, dt):
2     global brain
3     global reward
4     global conteo
5     global last_distance
6     global goal_x
7     global goal_y
8     global longueur
9     global largeur
10
11    longueur = self.width
12    largeur = self.height
13    if first_update:
14        init()
15
16    update_sand = True
17    if update_sand:
18        sand = np.zeros((int(self.width), int(
19        self.height)))
20        update_sand = False
21
22    # Diferencia de x-coordenadas entre la meta
23    # y el carrito.
24    xx = goal_x - self.car.x
25    # Diferencia de y-coordenadas entre la meta
26    # y el carrito.
27    yy = goal_y - self.car.y
28    orientation = Vector(*self.car.velocity).
29    angle((xx, yy))/180.

```

```

26     state = [orientation, self.car.signal1,
27               self.car.signal2, self.car.signal3]
28     action = brain.update(state, reward)
29
30     rotation = action2rotation[action]
31
32     self.car.move(rotation)
33
34     distance = np.sqrt((self.car.x - goal_x)**2
35 + (self.car.y - goal_y)**2)
36
37     # Actualizamos las posiciones de los
38     # sensores en el mapa.
39     self.ball1.pos = self.car.sensor1
40     self.ball2.pos = self.car.sensor2
41     self.ball3.pos = self.car.sensor3
42
43     self.scores.append(brain.score())
44
45     ***** ZONA DE RECOMPENSAS *****
46
47     # Cuando el carrito 1 se encuentra en
48     # obstaculo disminuye su velocidad a 1.
49     if sand[int(self.car.x), int(self.car.y)] >
50 0 or self.car.collide_widget(self.car2) or self.
51 car.collide_widget(self.car3) or self.car.
52 collide_widget(self.car4):
53     self.car.velocity = Vector(1, 0).rotate(
54 self.car.angle)
55     reward = -1
56
57     # Caso contrario el carrito1 mantiene una
58     # velocidad de 6.
59     else:
60         self.car.velocity = Vector(6, 0).rotate(
61 self.car.angle)
62         reward = -0.2
63         if distance < last_distance:
64             reward = 0.1
65
66     # Actualizamos lectura de sensores del carro
67     # 1 para que detecte cuando choca con otro carro.
68     if self.ball1.collide_widget(self.car2) or
69 self.ball1.collide_widget(self.car3) or self.
70 ball1.collide_widget(self.car4):
71         self.car.signal1 = 1.0
72
73     if self.ball2.collide_widget(self.car2) or
74 self.ball2.collide_widget(self.car3) or self.
75 ball2.collide_widget(self.car4):
76         self.car.signal2 = 1.0
77
78     if self.ball3.collide_widget(self.car2) or
79 self.ball3.collide_widget(self.car3) or self.
80 ball3.collide_widget(self.car4):
81         self.car.signal3 = 1.0
82
83     # Si el carrito logra salirse de los bordes del
84     # mapa, se asigna una penalidad de -1.
85     if self.car.x < 10:
86         self.car.x = 10
87         reward = -1
88     if self.car.x > self.width - 10:
89         self.car.x = self.width - 10
90         reward = -1
91     if self.car.y < 10:
92         self.car.y = 10
93         reward = -1
94     if self.car.y > self.height - 10:
95         self.car.y = self.height - 10
96         reward = -1
97
98     if distance < 100:
99         goal_x = self.width-goal_x

```

```
goal_y = self.height-goal_y
```

```
last_distance = distance
```

class Network(nn.Module):

Construcción de Red Neuronal Predictora de Q-values

Nuestra DQN tiene como eje principal una red neuronal que recibe un estado dado por la Orientación y las lecturas de los sensores y obtiene como salida los Q-values de las acciones disponibles para ese estado (Veáse figura 2), a continuación se muestra la implementación de nuestra red.

```

1 class Network(nn.Module):
2
3     def __init__(self, input_size, nb_action):
4         super(Network, self).__init__()
5         self.input_size = input_size
6         self.nb_action = nb_action
7         # Combinacion Lineal de toda la data
8         self.fc1 = nn.Linear(input_size, 30)
9         self.fc2 = nn.Linear(30, nb_action)
10
11     def forward(self, state):
12         x = F.tanh(self.fc1(state))
13         q_values = self.fc2(x)
14         return q_values

```

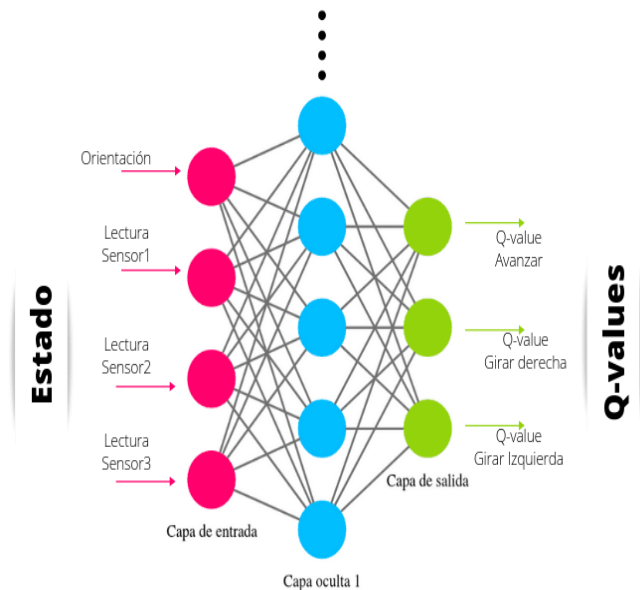


Figura 2. Representación de la Red Neuronal predictora de Q-values (Creación propia). Se muestra como un estado (Orientación, Sensor1, Sensor2, Sensor3) es ingresado a la red y se obtiene los Q-values de cada opción permitida.

VI. EXPERIMENTOS Y RESULTADOS

Inicialmente solo contamos con un carro, y diseñamos un mapa que se puede ver en figura 3, para el entrenamiento del carro. Y en la figura 4 se puede observar las recompensas a la lo largo del entrenamiento.

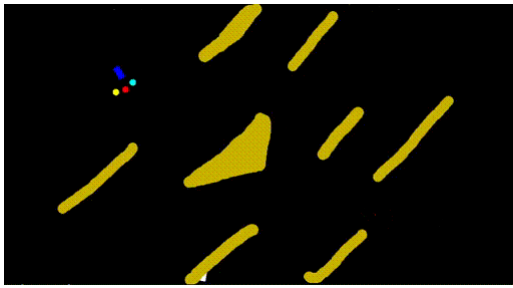


Figura 3. Mapa usado para el entrenamiento

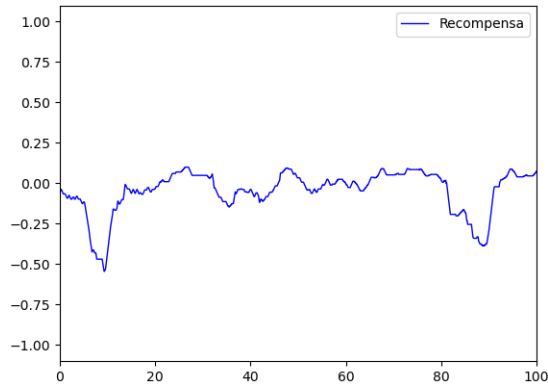


Figura 4. Recompensa a lo largo del entrenamiento

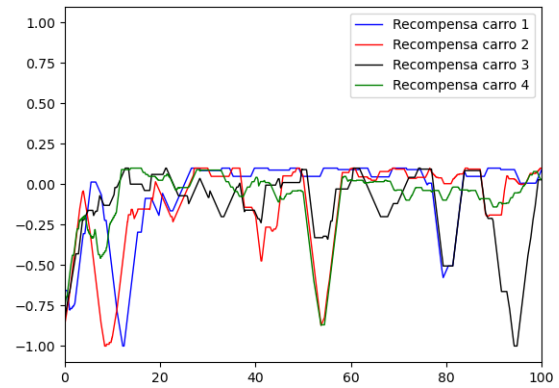


Figura 6. Recompensa a lo largo del entrenamiento

REFERENCIAS

- [1] Deep Q-Learning — An Introduction To Deep Reinforcement Learning. (2021). Retrieved 23 July 2021, from <https://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/>
- [2] Reinforcement Learning (DQN) Tutorial — PyTorch Tutorials 1.9.0+cu102 documentation. (2021). Retrieved 23 July 2021, from https://pytorch.org/tutorials/intermediate/reinforcement_q_learning.html
- [3] Welcome to Kivy — Kivy 2.0.0 documentation. (2021). Retrieved 25 July 2021, from <https://kivy.org/doc/stable/>

En nuestra implementación hemos agregado 3 carros más en el mismo mapa de entrenamiento como se puede ver en la figura 5, de la misma manera se midió las recompensas a lo largo del entrenamiento como se muestra en la figura 6.



Figura 5. Mapa usado para el entrenamiento

VII. CONCLUSIONES

- Deep Q-Learning no es perfecto para la solución de este problema, más tiene un increíble acercamiento al objetivo de un "Vehículo Autónomo"
- Las funciones de activación influyen fuertemente a la eficiencia del vehículo, mostrando mejores resultados con una función tanh que con una función RELu durante las pruebas.