

---

**Instituto Tecnológico de Costa Rica****Escuela de Ingeniería Electrónica****Trabajo Final de Graduación****Proyecto:** Método basado en aprendizaje reforzado para el control automático de una planta no lineal.**Estudiante:** Oscar Andrés Rojas Fonseca

I Semestre 2024

---

**Firma del asesor**

---

**Bitácora de trabajo**

Fecha	Actividad	Anotaciones	Horas dedicadas
15/04/2024	<b>1.</b> Redefinición de la conversión del código para valores discretos ( <i>CartPole</i> ) a valores continuos ( <i>Pendulum</i> ).	a) El error en <i>select_action()</i> se corrigió pero desconfiguró parte de la función <i>optimize_model()</i> . Corrección del error. b) Persisten los problemas de indexado y proceso.	6 horas
15/04/2024	<b>2.</b> Pruebas de entrenamiento del modelo ( <i>Pendulum</i> ).	a) Se entrenaron cuatro modelos diferentes a 600 episodios para comparar el efecto de cuatro propuestas de redes neuronales artificiales (ANN).	4 horas
16/04/2024	<b>3.</b> Búsqueda de la teoría de los métodos <i>PPO</i> y <i>actor – critic</i> dada la necesidad del manejo del <i>action space</i> con valores continuos.	a) Revisión general de los conceptos para ambos casos. Fuente importante en publicación de <i>PPO</i> [1]. a) La fuente [1] facilita un enlace al repositorio de GitHub respectivo con ejemplos debidamente documentados.	4 horas
17/04/2024	<b>4.</b> Reunión de seguimiento con el asesor del proyecto.	a) Revisión de avance en el código y errores de forma. b) Se acordó continuar con el interés en los métodos como <i>PPO</i> como opción.	2 horas

17/04/2024	<b>5.</b> Búsqueda de métodos para el manejo de valores continuos en <i>DQN</i> .	<p>a) Las opciones popularmente mencionadas son el <i>actor-critic</i> y el <i>DDPG</i>. Se procede a investigar <i>DDPG</i></p> <p>b) La opción de discretizar el <i>action space</i> también se menciona.</p>	4 horas
17/04/2024	<b>6.</b> Prueba de discretización del <i>action space</i> del env <i>Pendulum</i> .	<p>a) Se logró adaptar el código del <i>PendulumDQN</i> a una versión discretizada <i>PendulumDQN_discrete</i>, depende principalmente de la resolución seleccionada (<i>n_actions</i>).</p> <p>b) Pruebas de entrenamiento de hasta 100 episodios.</p>	6 horas
19/04/2024	<b>7.</b> Revisión de implementación del <i>DDPG</i> y el <i>PPO</i> .	<p>a) Se buscaron ejemplos para aplicar <i>DDPG</i> mediante PyTorch para evitar ajustes innecesarios.</p> <p>b) Estudio de otros repositorios para aplicar <i>PPO</i>.</p>	4 horas
20/04/2024	<b>8.</b> Montaje y primera prueba del código <i>PPO</i> para <i>Pendulum</i> .	<p>a) Revisión del error por cambio de versión <i>Gym</i> a <i>Gymnasium</i>.</p> <p>b) Implementación del procesamiento <i>CUDA</i> en el código.</p>	8 horas
21/04/2024	<b>9.</b> Prueba de entrenamiento con versión base del <i>PPO</i> para <i>Pendulum</i> .	<p>a) Entrenamiento del modelo con <i>render("human")</i> de 200k episodios. Mal desempeño.</p> <p>b) Entrenamiento del modelo con <i>render("rgb_array")</i> de 200M de episodios. Interrumpido por tiempo.</p>	6 horas

21/04/2024	<b>10.</b> Adición de la función <i>calculate_reward()</i> al código base <i>PPO</i> de <i>Pendulum</i> .	a) Ajuste de la función para retornar valores <i>float</i> , no <i>torch.tensor()</i> . b) Ajuste de las funciones <i>log_summary()</i> para una mejor apreciación del proceso de entrenamiento y prueba. c) Estudio del uso de <i>Tensorboard</i> para procesos de entrenamiento de <i>Machine Learning</i> .	6 horas
22/04/2024	<b>11.</b> Pruebas de implementación del método <i>DDPG</i> para <i>Pendulum</i> .	a) Montaje del código en carpeta <i>PruebasDDPG</i> y se realizaron las primeras pruebas de entrenamiento. Errores de forma en el proceso.	4 horas
22/04/2024	<b>12.</b> Pruebas de entrenamiento de <i>PendulumPPO</i> .	a) Entrenamiento de modelo con <i>1M</i> de episodios. Se muestra con avance prometedore en cuanto a las <i>rewards</i> .	6 horas
23/04/2024	<b>13.</b> Pruebas montaje e incorporación del código <i>pahm_model.py</i> proporcionado por el profesor asesor.	a) Primeras pruebas de implementación y adición del proceso método <i>PPO</i> al entrenamiento.	4 horas
Total de horas de trabajo:			64 horas

## Contenidos de actividades

### Primera prueba de entrenamiento *PPO*

El código base utilizado [1] contaba con un valor de *timesteps* por defecto de docientos millones (*200M*), por lo que luego de adaptar el código para su funcionamiento mediante el procesamiento *CUDA*, se probó cumplir con esta cantidad, como se observa en la Figura 1.

Sin embargo y apesar de observarse un avance en el entrenamiento, luego de un aproximado de tiempo de cuatro horas transcurridas se proyectó un faltante importante, por lo que se decidió interrumpir el proceso para lograr mejoras en otros puntos del proyecto.

```
main(args)
[5] 150m 29.7s
... iteration LOOK: 0.17 SECS
-----

----- Iteration #1220 -----
Average Episodic Length: 200.0
Average Episodic Return: -145.06
Average Loss: -0.00434
Timesteps So Far: 2684000
Iteration took: 6.98 secs
-----

----- Iteration #1221 -----
Average Episodic Length: 200.0
Average Episodic Return: -157.11
Average Loss: -0.00127
Timesteps So Far: 2686200
Iteration took: 8.02 secs
-----
```

Figure 1: Proceso de entrenamiento del modelo para *PendulumPPO* con 200M de episodios.

## Ajuste en la función *log\_summary()* de *PPO*

Se buscó una mejora en la forma de monitoreo del avance del entrenamiento del modelo, por lo que agregó la lógica de la librería *Matplotlib* y se probó con un entrenamiento inicial de 100k episodios como se muestra en la Figura 2.

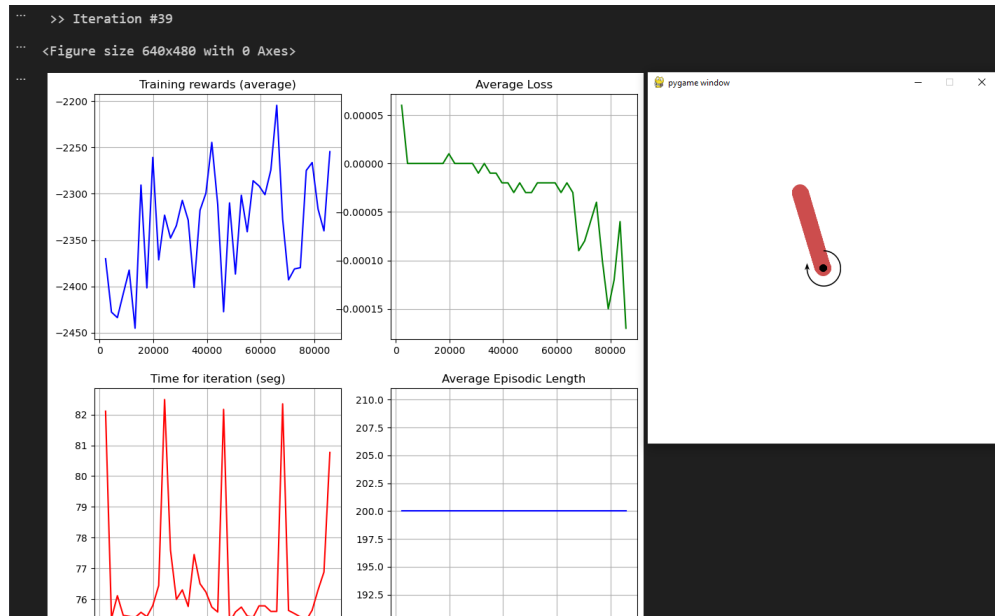


Figure 2: Proceso de entrenamiento del modelo para *PendulumPPO* con 100k episodios.

También se empezó con pruebas simples de monitoreo mediante la librería *Tensorboard*, las cuales se encuentran en proceso.

## Entrenamiento de modelo *PPO* de 1M de episodios

Luego de ajustar el monitoreo y agregar la lógica del *target\_angle* al código, se procedió a entrenar un modelo con 1 millón de episodios y fijar el ángulo en  $135^\circ$ , Figura 3, mostrando un desempeño prometedor.

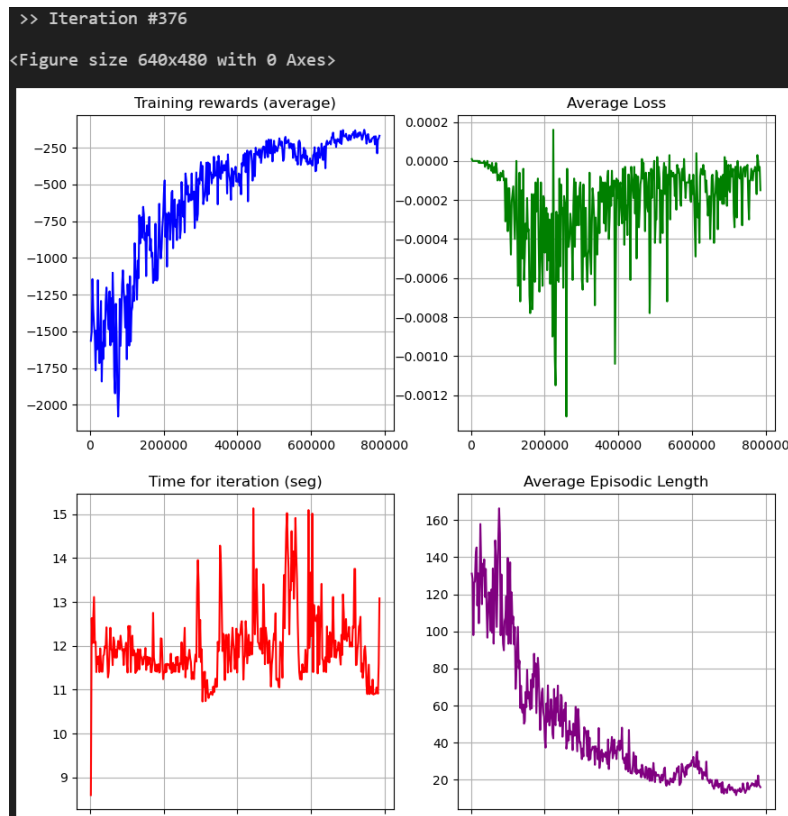


Figure 3: Proceso de entrenamiento del modelo para *PendulumPPO* con 1M de episodios y *target\_angle* = 135.

Se encuentra en proceso de implementación la variación constante del ángulo y más entrenamientos, dada la necesidad de periodos largos de tiempo para entrenamientos de modelos prometedores.

## Referencias

- [1] E. Yang-Yu, “Coding ppo from scratch with pytorch (part 1/4),” *Medium*, 2020.