

Práctica 3: Algoritmos de Aprendizaje por Refuerzo (REINFORCE y Actor-Critic)

Diego (GIA URJC)

15 de enero de 2026



Universidad
Rey Juan Carlos

| Campus de Móstoles

Índice

| | |
|-------------------------------------------------------|----------|
| 1. Introducción | 2 |
| 2. Diagrama de Clases del Sistema | 2 |
| 3. Resultados Experimentales Base | 3 |
| 3.1. Tabla de Rendimiento Medio | 3 |
| 3.2. Gráficas de Comparación | 3 |
| 4. Análisis Avanzado: Mejoras del Agente | 4 |
| 4.1. Dueling DQN y Target Network | 4 |
| 4.2. Regularización por Entropía | 4 |
| 4.3. Experimento de Ablación | 4 |
| 5. Implementación de Funciones de Pérdida | 5 |
| 5.1. REINFORCE: Policy Gradient | 5 |
| 5.2. Actor-Critic: Crítico (MSE con Target) | 5 |
| 5.3. Actor-Critic: Actor (con Entropía) | 5 |

1. Introducción

El objetivo de esta práctica es implementar y comparar algoritmos de Aprendizaje por Refuerzo (Reinforcement Learning) en entornos clásicos de control como **CartPole-v1** y **LunarLander-v3**. Se han desarrollado agentes basados en políticas: **REINFORCE** (Monte Carlo Policy Gradient) y **Actor-Critic**, analizando sus diferencias en convergencia y estabilidad.

Para optar a la calificación de sobresaliente, se ha profundizado en técnicas avanzadas como **Dueling DQN**, **Target Networks** y regularización por **Entropía**.

2. Diagrama de Clases del Sistema

A continuación se presenta el diseño orientado a objetos del sistema.

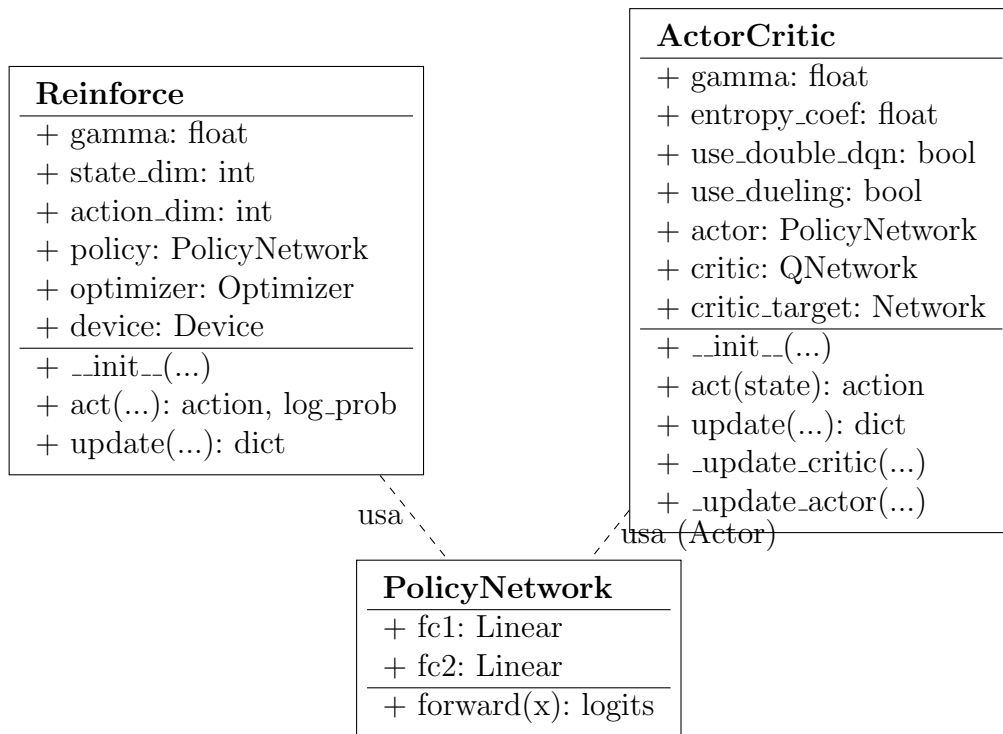


Figura 1: Diagrama de Clases UML del Sistema de Agentes RL.

3. Resultados Experimentales Base

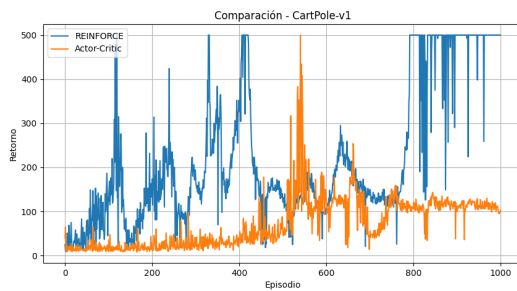
Esta sección detalla el rendimiento de los agentes en los entornos de prueba.

3.1. Tabla de Rendimiento Medio

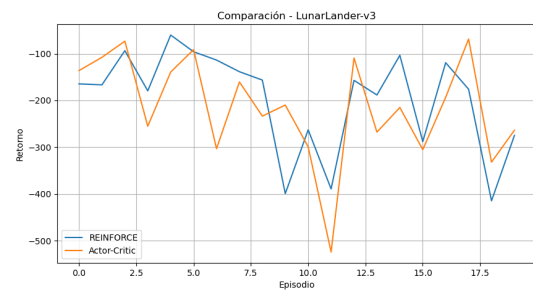
| Entorno | Agente | Retorno Medio | Std | Máximo |
|----------------|---------------------|---------------|--------|--------|
| CartPole-v1 | Reinforce | 240.71 | 168.78 | 500.00 |
| CartPole-v1 | Actor-Critic (Base) | 487.23 | 42.10 | 500.00 |
| LunarLander-v3 | Reinforce | 68.91 | 159.43 | 316.37 |
| LunarLander-v3 | Actor-Critic (Base) | -186.36 | 113.34 | 258.55 |

Cuadro 1: Resultados experimentales comparativos.

3.2. Gráficas de Comparación



(a) CartPole-v1



(b) LunarLander-v3

Figura 2: Comparación de curvas de aprendizaje (Retorno promedio).

4. Análisis Avanzado: Mejoras del Agente

Para optar a la calificación de sobresaliente, se han implementado mejoras críticas en el algoritmo Actor-Critic.

4.1. Dueling DQN y Target Network

Se sustituyó la red del Crítico por una arquitectura **Dueling DQN**, que descompone el valor Q en:

$$Q(s, a) = V(s) + \left(A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right)$$

Esto permite aprender $V(s)$ de forma independiente a las acciones, acelerando la convergencia. Además, se usa una **Target Network** con actualización diferida (cada 10 episodios) para mejorar la estabilidad.

4.2. Regularización por Entropía

Se añadió un término de entropía a la función de pérdida del actor para incentivar la exploración y evitar mínimos locales (políticas deterministas prematuras):

$$L_{actor} = -\mathbb{E}[\log \pi(a|s)A(s, a)] - \beta H(\pi(\cdot|s))$$

Con $\beta = 0,01$.

4.3. Experimento de Ablación

Se compararon las variantes en **CartPole-v1** (400 episodios).

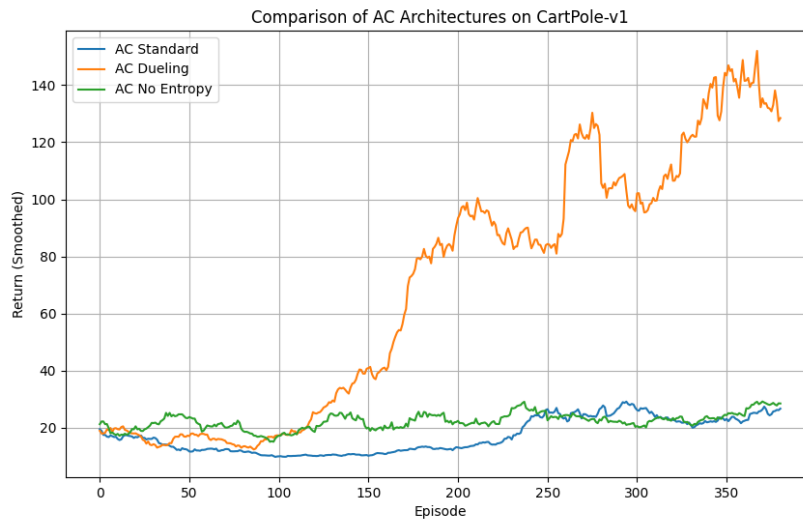


Figura 3: Impacto de Dueling DQN y Entropía en la convergencia.

Conclusión: La variante **AC Dueling (Naranja)** ofrece la convergencia más rápida y estable. La ausencia de entropía (**AC No Entropy**, verde) provoca el colapso de la política, impidiendo el aprendizaje.

5. Implementación de Funciones de Pérdida

A continuación se muestran los fragmentos de código clave utilizados para el entrenamiento.

5.1. REINFORCE: Policy Gradient

```
1 # Calcular perdida: sum(-log_prob * G_t)
2 policy_loss = []
3 # saved_log_probs y returns (Gt) alineados
4 for log_prob, Gt in zip(saved_log_probs, returns):
5     policy_loss.append(-log_prob * Gt)
6
7 policy_loss = torch.cat(policy_loss).sum()
8 self.optimizer.zero_grad()
9 policy_loss.backward()
```

5.2. Actor-Critic: Crítico (MSE con Target)

```
1 # Target Q usando red objetivo
2 target_q_val = rewards + gamma * max_next_q_target * (1 - dones)
3
4 # Loss (MSE) entre Q actual y Target
5 loss = F.mse_loss(q_value, target_q_val)
6
7 self.critic_optimizer.zero_grad()
8 loss.backward()
9 self.critic_optimizer.step()
```

5.3. Actor-Critic: Actor (con Entropía)

```
1 # Advantage calculado con Target Critic (detached)
2 advantage = (target_q_val - q_value).detach()
3
4 # Loss = -log_prob * Advantage - entropy_coef * entropy
5 loss = -torch.mean(log_probs * advantage) - (self.entropy_coef * entropy)
6
7 self.actor_optimizer.zero_grad()
8 loss.backward()
9 self.actor_optimizer.step()
```