

Aprendizaje por Refuerzo

Materia: Redes Neuronales.

Abril 2020

Bibliografía

- › Hands-on ML with Sk-learn & TensorFlow. (Géron)
- › Reinforcement Learning: an introduction. (Sutton & Barto)
- › Foundations of ML. (Mohri)
- › Algorithms for RL. (Szepesvari)
- › <https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf>
- › Link:
<https://drive.google.com/drive/folders/1HD-ty2o09S4SpappK4LVK8rPdpxA5Zud?usp=sharing>

Aprendizaje por refuerzo (RL)

1. Campo de Machine Learning (ML) – 1950

Aprendizaje por refuerzo (RL)

1. Campo de Machine Learning (ML) – 1950
2. Objetivo: Mapear situaciones a acciones (¿Qué debo hacer?) para maximizar una recompensa.

Aprendizaje por refuerzo (RL)

1. Campo de Machine Learning (ML) – 1950
2. Objetivo: Mapear situaciones a acciones (¿Qué debo hacer?) para maximizar una recompensa.
3. Estudio de planificación y aprendizaje en un escenario.

Aprendizaje por refuerzo (RL)

1. Campo de Machine Learning (ML) – 1950
2. Objetivo: Mapear situaciones a acciones (¿Qué debo hacer?) para maximizar una recompensa.
3. Estudio de planificación y aprendizaje en un escenario.
4. Elementos: Un aprendiz (**agente**) interactúa con un **ambiente** para llegar a un objetivo medido en **recompensa**.

Aprendizaje por refuerzo (RL)

1. Campo de Machine Learning (ML) – 1950
2. Objetivo: Mapear situaciones a acciones (¿Qué debo hacer?) para maximizar una recompensa.
3. Estudio de planificación y aprendizaje en un escenario.
4. Elementos: Un aprendiz (**agente**) interactúa con un **ambiente** para llegar a un objetivo medido en **recompensa**.
5. Es un enfoque distinto a los aprendizajes supervisado y no supervisados (SL, UL).
Se recibe un data set pasivamente desde un supervisor externo. **La interacción provee datos activamente.**

SL: (Descripción X, Etiqueta Y) \Rightarrow Extrapolar.

UL: (Datos X) \Rightarrow Estructura oculta.

RL: Interacción \Rightarrow Maximizar una recompensa.

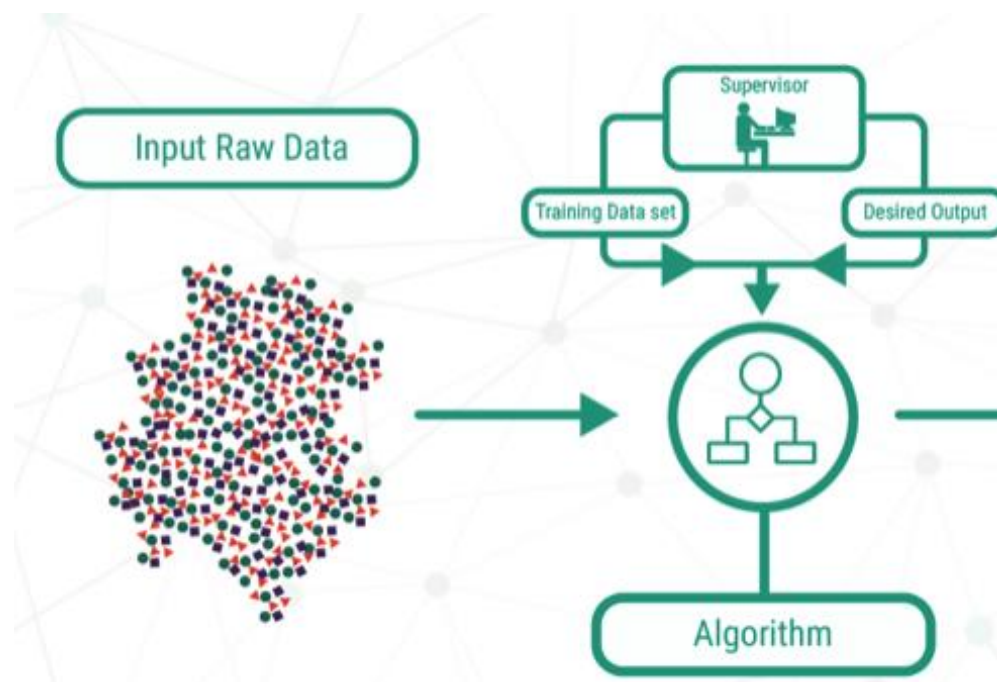
No hay representantes a todas las situaciones. Es necesaria, la experiencia.

Aprendizaje por refuerzo (RL)

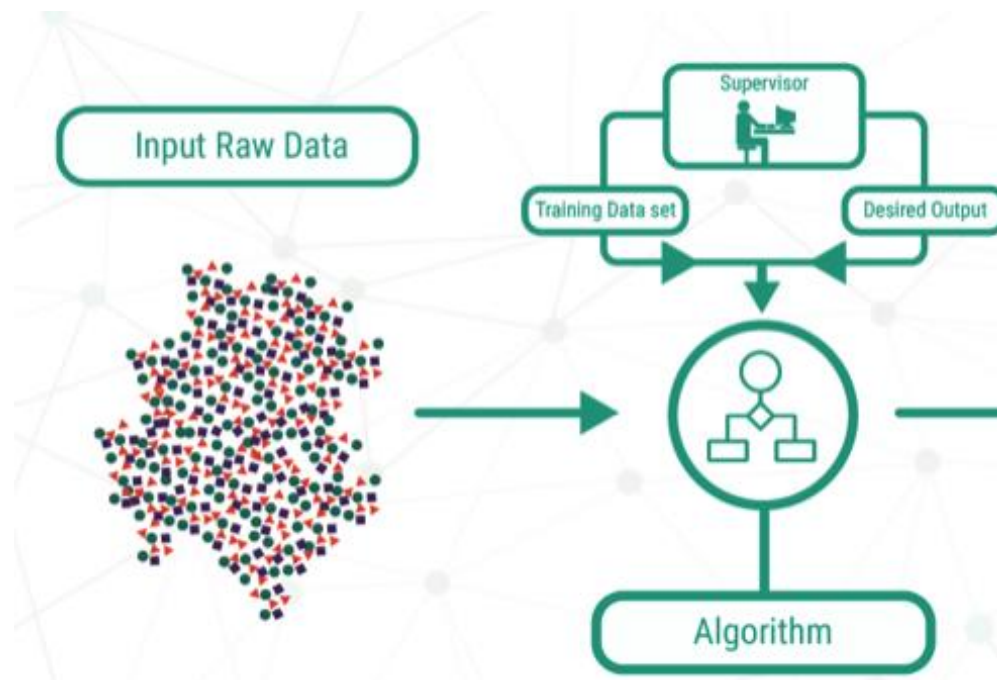
1. Campo de Machine Learning (ML) – 1950
2. Objetivo: Mapear situaciones a acciones (¿Qué debo hacer?) para maximizar una recompensa.
3. Estudio de planificación y aprendizaje en un escenario.
4. Elementos: Un aprendiz (**agente**) interactúa con un **ambiente** para llegar a un objetivo medido en **recompensa**.
5. Es un enfoque distinto a los aprendizajes supervisado y no supervisados (SL, UL).
Se recibe un data set pasivamente desde un supervisor externo. **La interacción provee datos activamente.**

SL: (Descripción X, Etiqueta Y) \Rightarrow Extrapolar.
UL: (Datos X) \Rightarrow Estructura oculta.
RL: Interacción \Rightarrow Maximizar una recompensa.
No hay representantes a todas las situaciones. Es necesaria, la experiencia.
6. Aplicaciones en teoría de control, optimización, ciencia cognitivas. Games: Ataris, Alpha-Go

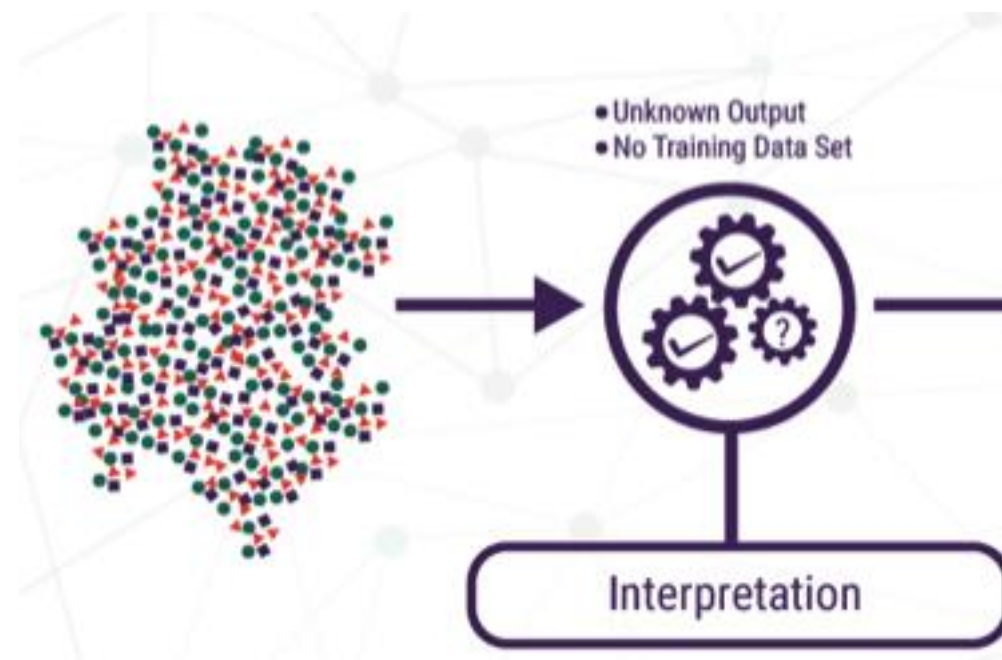
SUPERVISED LEARNING



SUPERVISED LEARNING

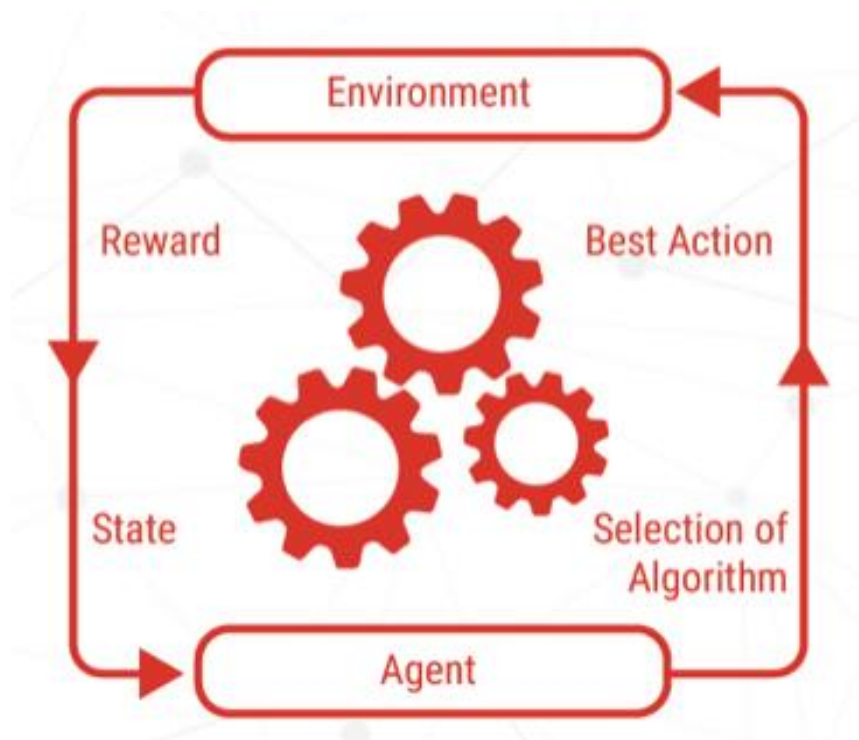


UNSUPERVISED LEARNING



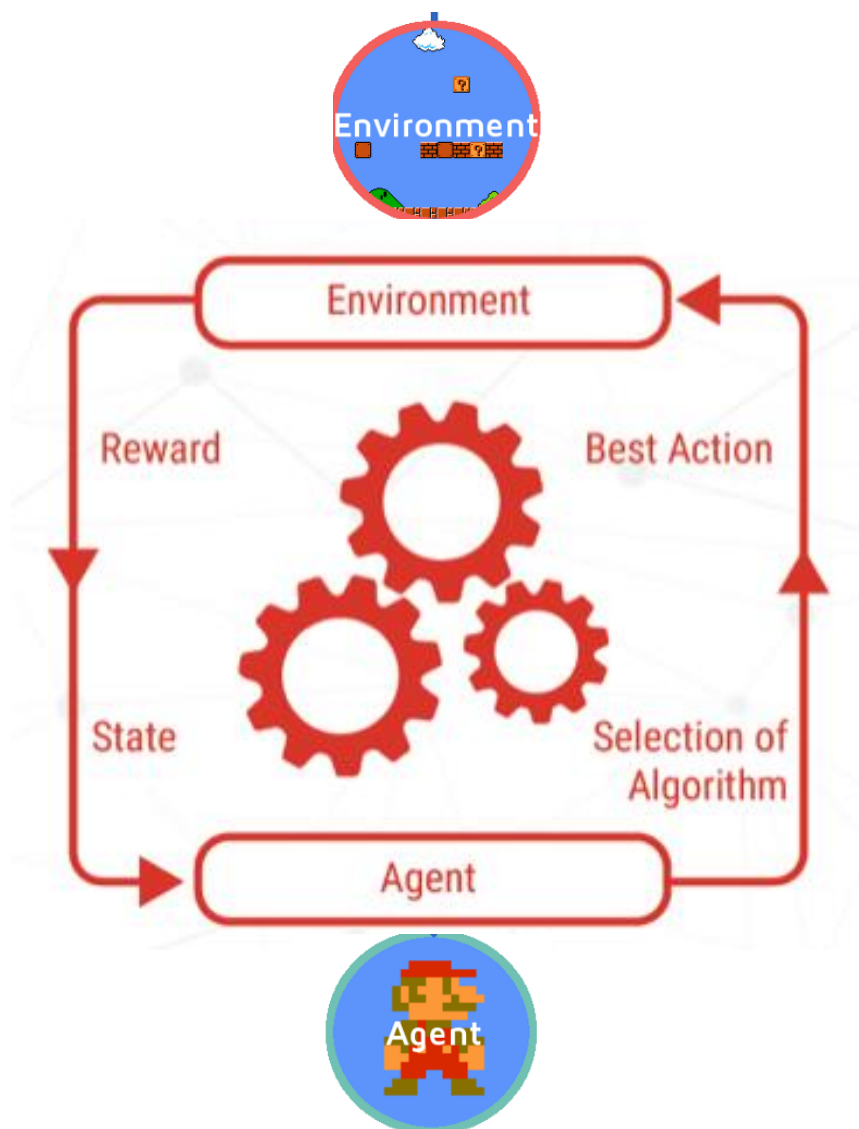
π

REINFORCEMENT LEARNING



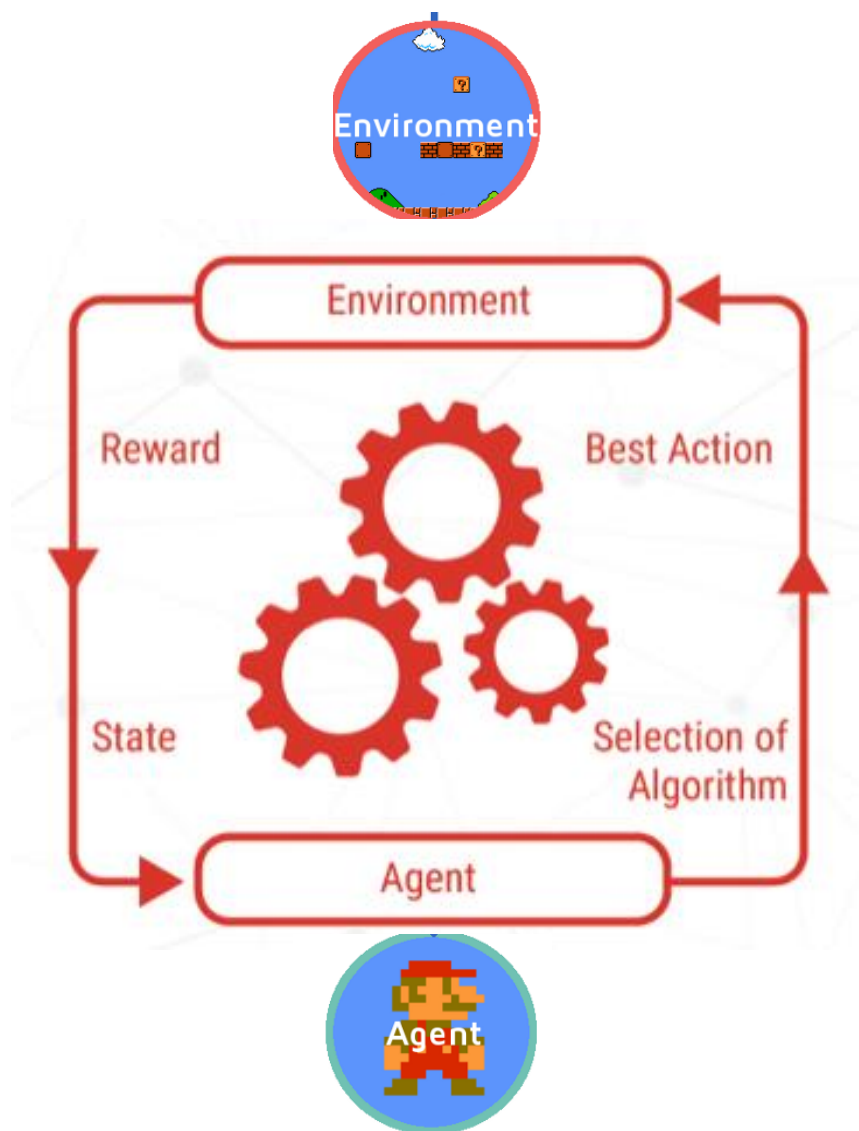
π

REINFORCEMENT LEARNING



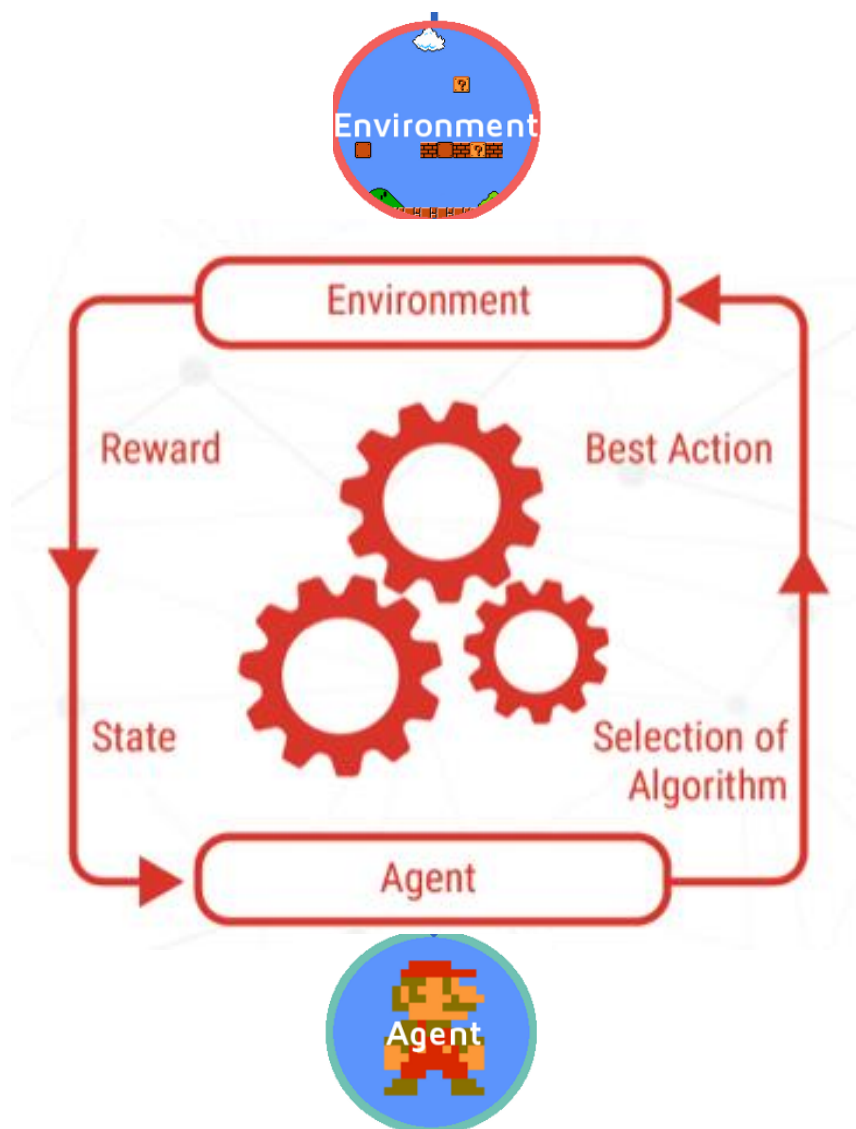
π

REINFORCEMENT LEARNING



Agente:

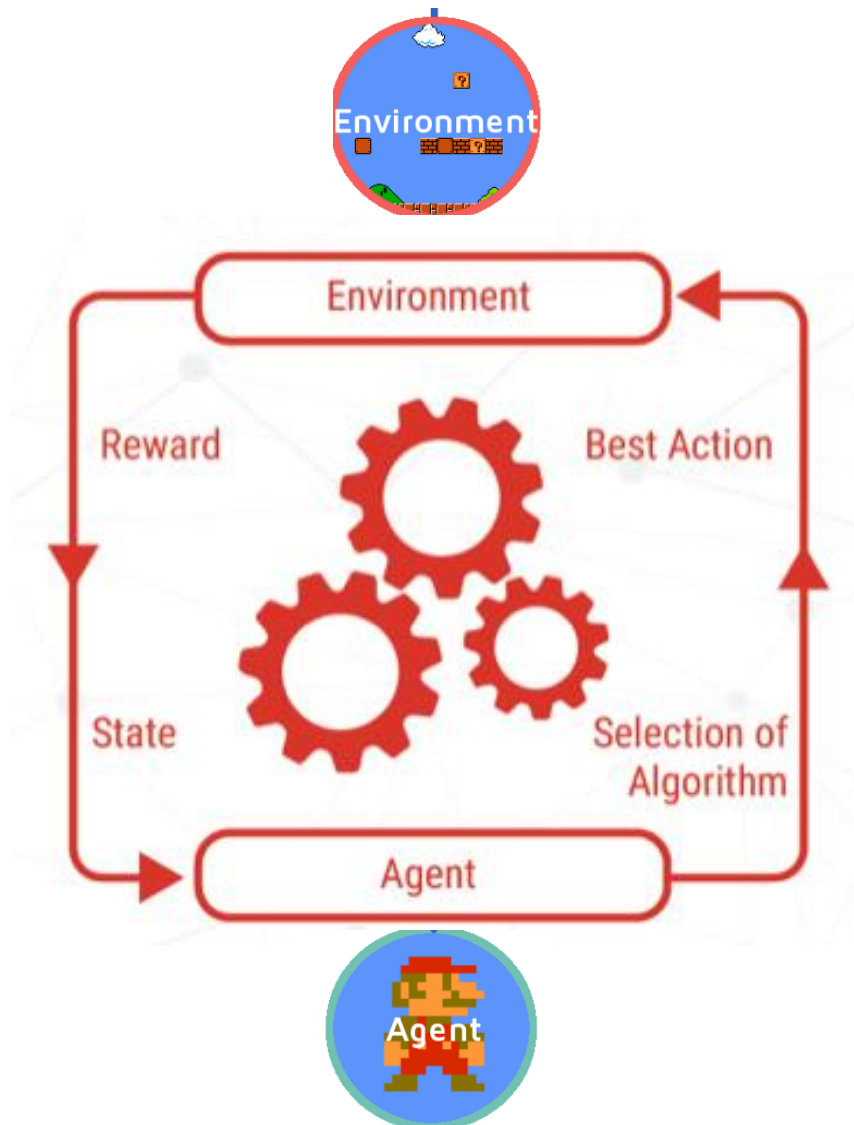
REINFORCEMENT LEARNING



Agente: Controlador del personaje Mario.

Ambiente:

REINFORCEMENT LEARNING

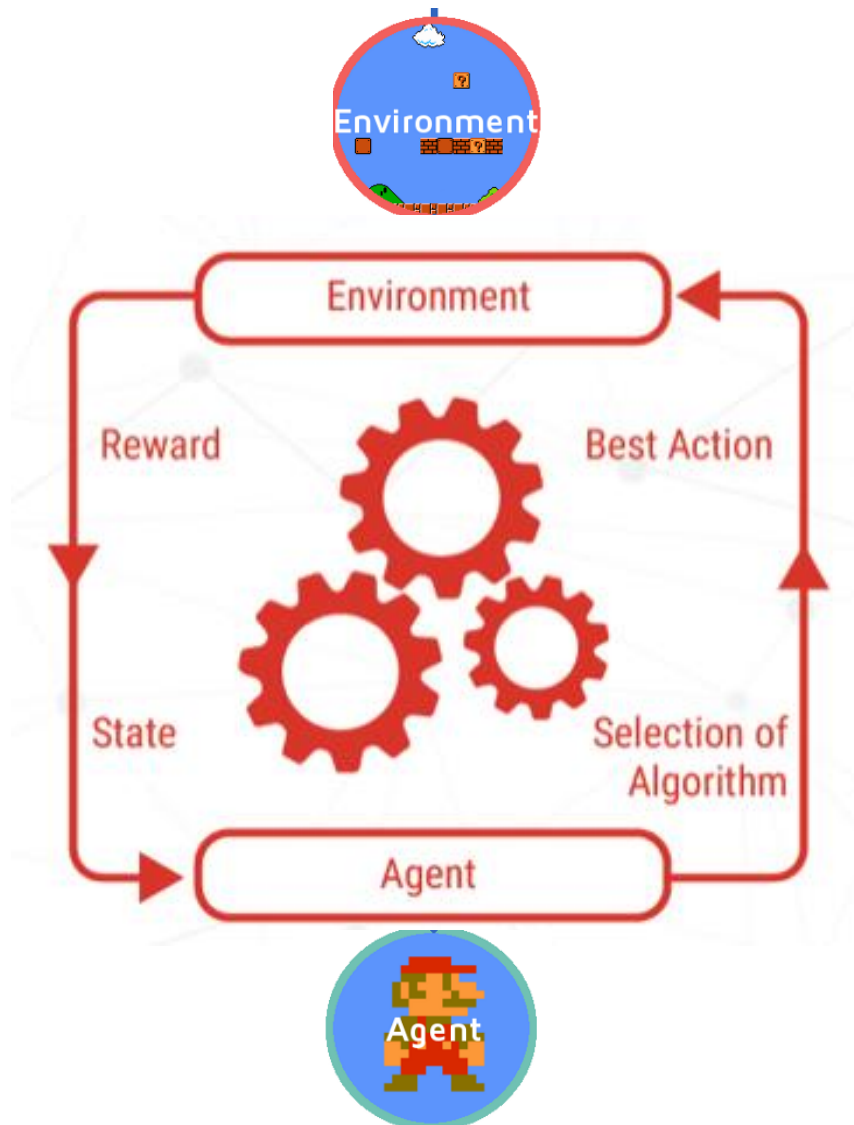


Agente: Controlador del personaje Mario.

Ambiente: Simulación del entorno del juego Mario Bros.

Observación:

REINFORCEMENT LEARNING



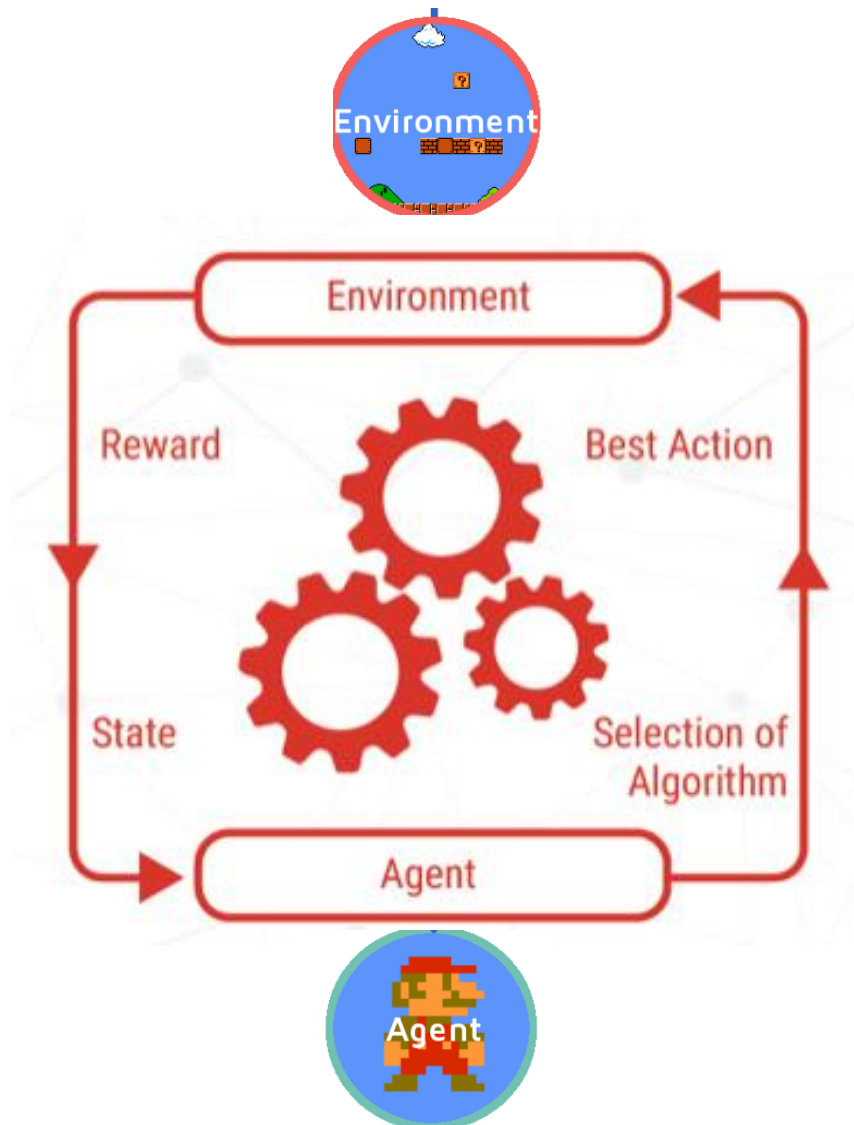
Agente: Controlador del personaje Mario.

Ambiente: Simulación del entorno del juego Mario Bros.

Observación: Screenshot – Imagen

Acción:

REINFORCEMENT LEARNING



Agente: Controlador del personaje Mario.

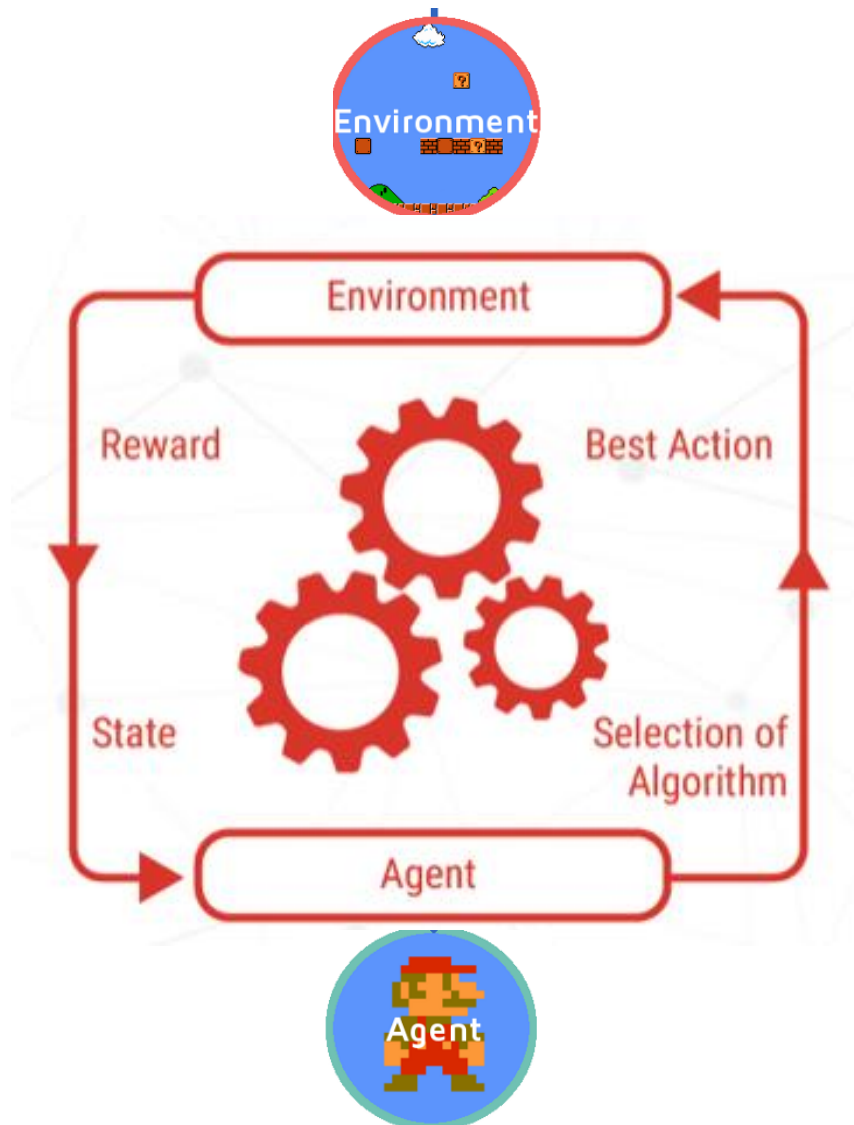
Ambiente: Simulación del entorno del juego Mario Bros.

Observación: Screenshot – Imagen

Acción: Posiciones del Joystick (Δ ∇ \triangleright \triangleleft)

Recompensa:

REINFORCEMENT LEARNING



Agente: Controlador del personaje Mario.

Ambiente: Simulación del entorno del juego Mario Bros.

Observación: Screenshot – Imagen

Acción: Posiciones del Joystick (Δ ∇ \triangleright \triangleleft)

Recompensa: Puntos del juego (Cantidad de vidas, Monedas)

π

Política

π

Política

Es el algoritmo usado por el agente para las acciones a tomar.

Mapeo entre observación y la acción elegida.

Política

 π

Es el algoritmo usado por el agente para las acciones a tomar.

Mapeo entre observación y la acción elegida.

PROBLEMA 1:

Un programa controla la caminata de un robot (agente) y observa su alrededor (ambiente) mediante cámaras (observación). Las acciones tomadas por el agente son señales que activan el sistema motor del robot.

π

Política

Es el algoritmo usado por el agente para las acciones a tomar.

Mapeo entre observación y la acción elegida.

PROBLEMA 1:

Un programa controla la caminata de un robot (agente) y observa su alrededor (ambiente) mediante cámaras (observación). Las acciones tomadas por el agente son señales que activan el sistema motor del robot.

Ejemplos de Políticas:

Determinista, sin considerar la observación:

Determinista, considerando la observación:

Estocástica, sin considerar la observación:

Estocástica, considerando la observación:

π

Política

Es el algoritmo usado por el agente para las acciones a tomar.

Mapeo entre observación y la acción elegida.

PROBLEMA 1:

Un programa controla la caminata de un robot (agente) y observa su alrededor (ambiente) mediante cámaras (observación). Las acciones tomadas por el agente son señales que activan el sistema motor del robot.

Ejemplos de Políticas:

Determinista, sin considerar la observación:

A todo tiempo, el robot camina derecho (no dobla).

Determinista, considerando la observación:

El robot caminará derecho si no hay obstáculos al frente. En caso contrario, girará 90° a la izquierda.

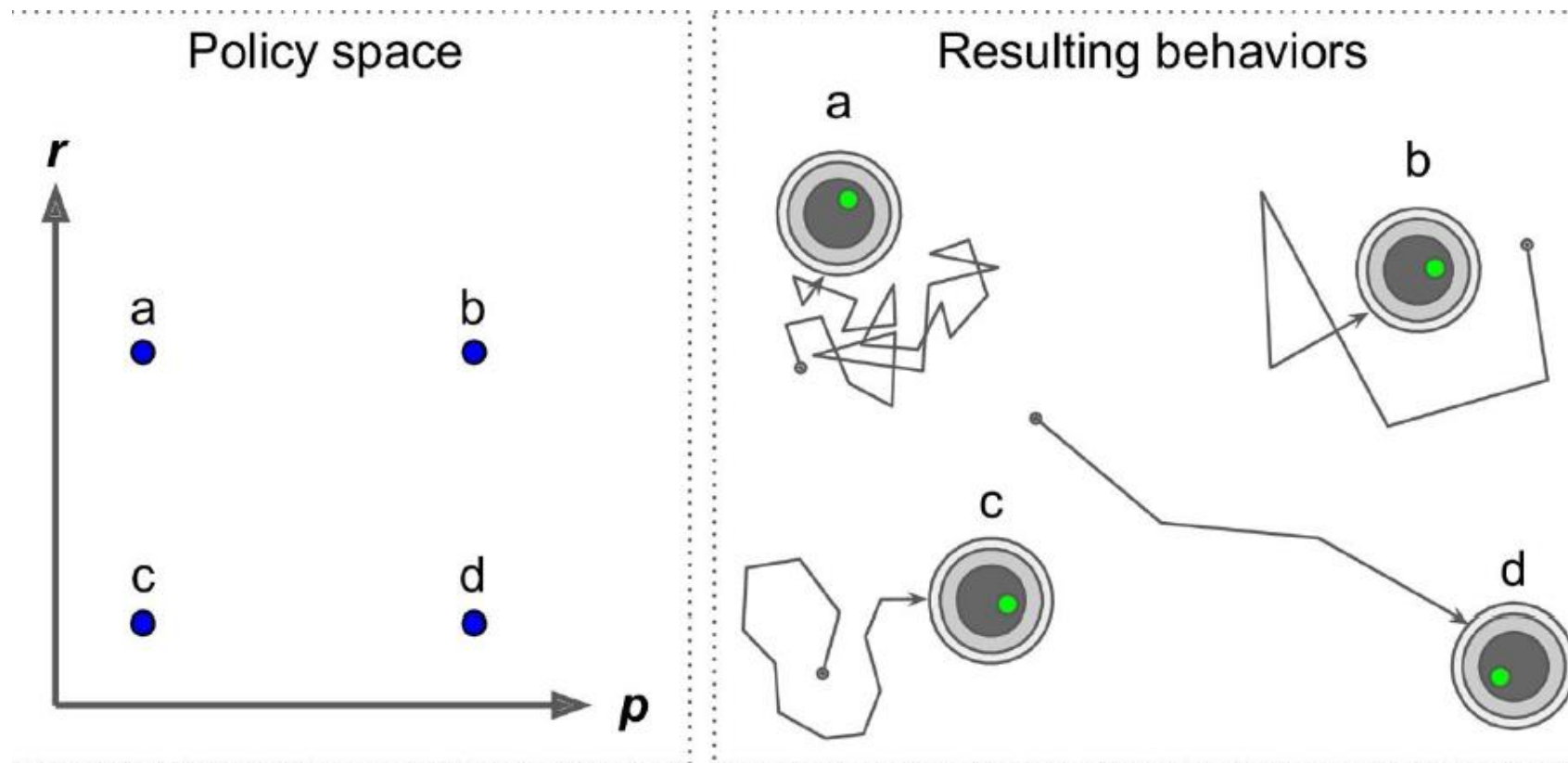
Estocástica, sin considerar la observación:

El robot caminará derecho con probabilidad p ó girará $r \in [\pi, -\pi]$ con probabilidad $1-p$

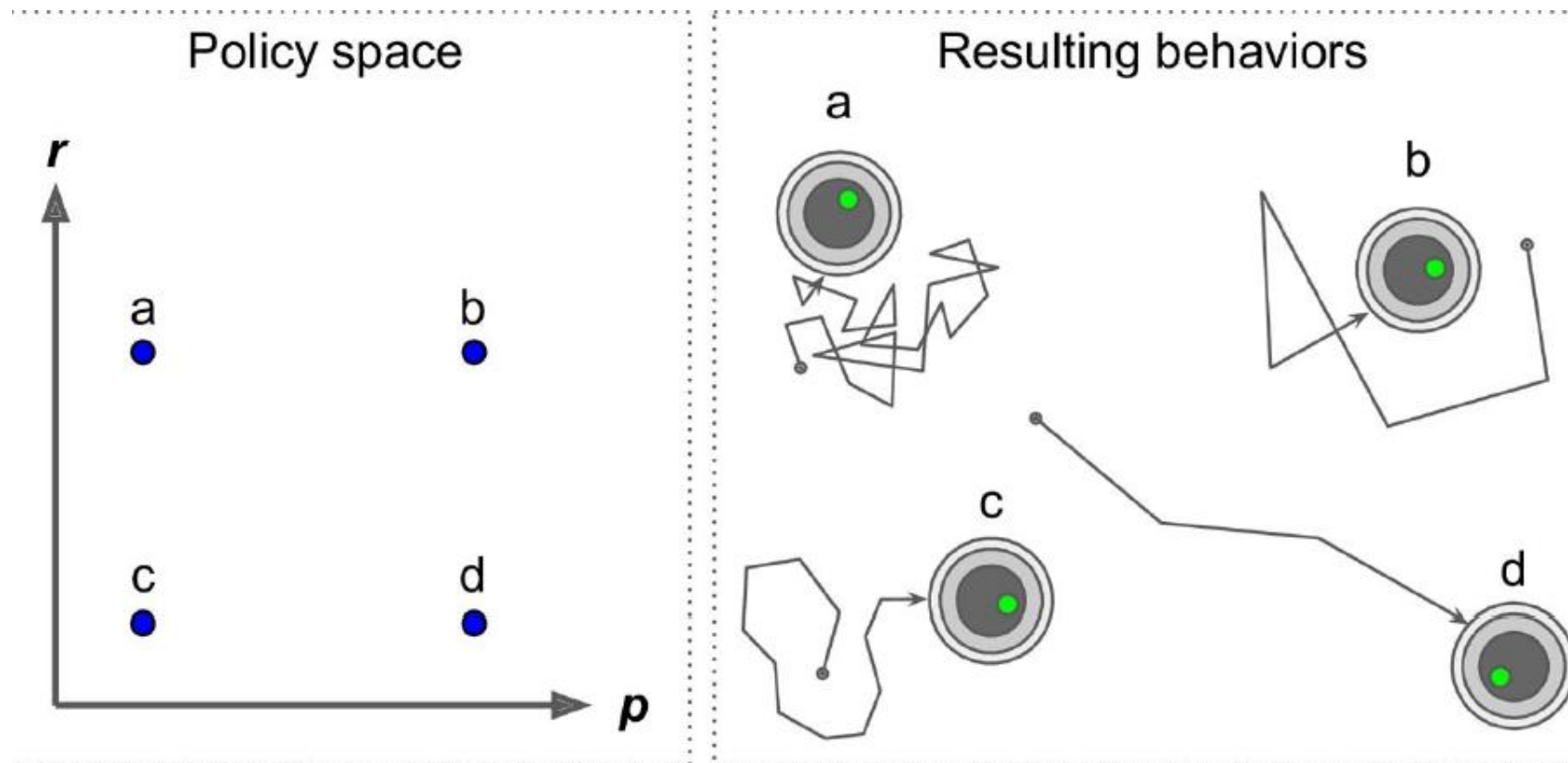
Estocástica, considerando la observación:

El robot caminará derecho si no hay obstáculos al frente. En otro caso, girará un ángulo random.

Política

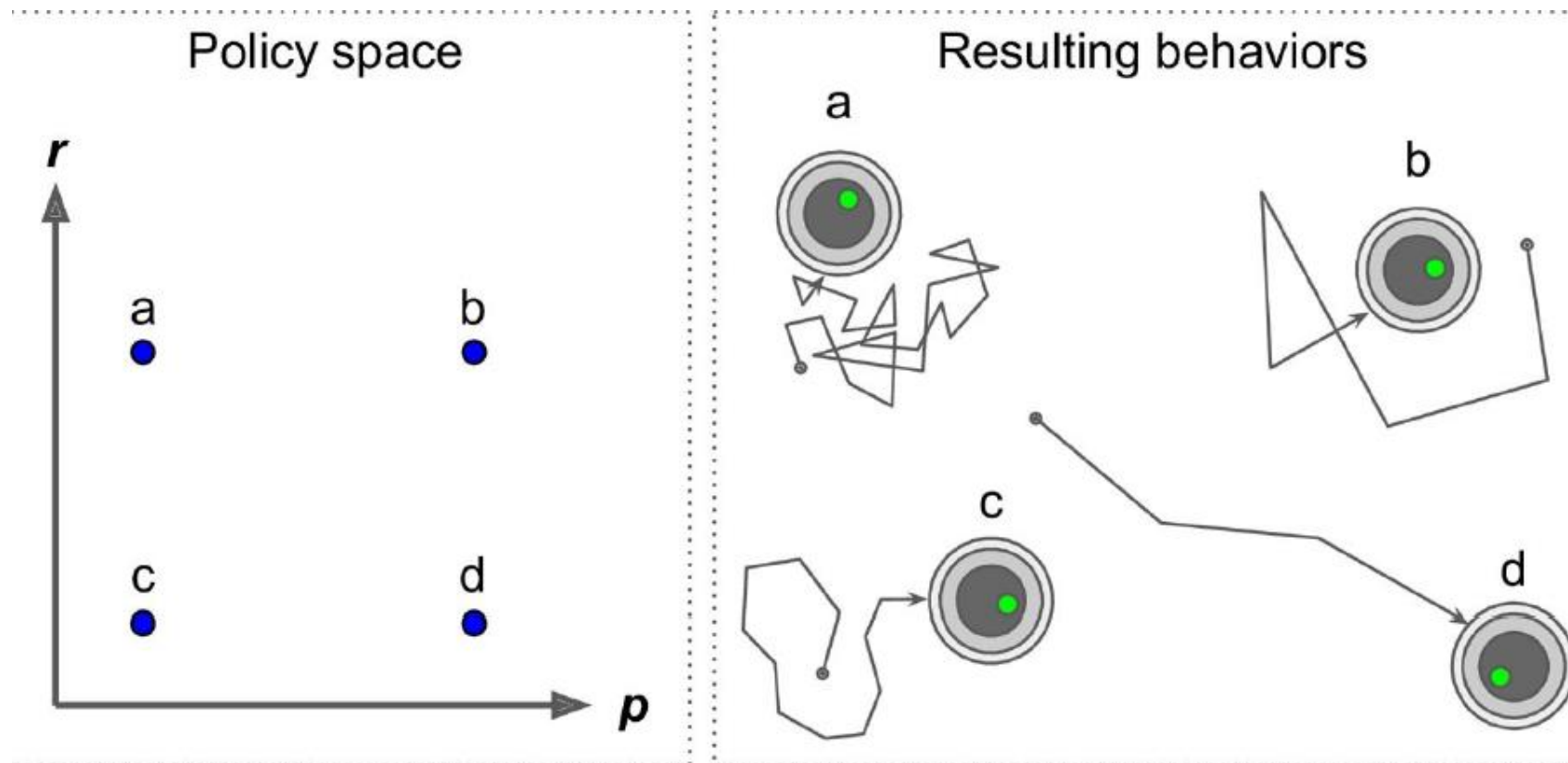


Política



PROBLEMA 2: Si la recompensa disminuye con el tiempo de llegada a un punto fijo desde otro dado, ¿Cuáles son los valores de (r,p) óptimos?

Política



PROBLEMA 2: Si la recompensa disminuye con el tiempo de llegada a un punto fijo desde otro dado, ¿Cuáles son los valores de (r,p) óptimos?

$$(r,p)^* = \operatorname{argmax}_{(r,p)} R = \operatorname{argmin}_{(r,p)} T$$

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

PROBLEMA 3:

* ¿Cómo se resuelve analíticamente este problema?

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

PROBLEMA 3:

* ¿Cómo se resuelve analíticamente este problema?

La búsqueda de máximos de funciones multivariantes se basa en resolver el gradiente nulo y hessiana definida negativa.

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

PROBLEMA 3:

* ¿Cómo se resuelve analíticamente este problema?

La búsqueda de máximos de funciones multivariantes se basa en resolver el gradiente nulo y hessiana definida negativa.

* ¿Qué dificultad se presenta con el aumento de la dimensiones N de parámetros de la función de recompensa? ¿Qué métodos permiten solucionarlo?

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

PROBLEMA 3:

* ¿Cómo se resuelve analíticamente este problema?

La búsqueda de máximos de funciones multivariantes se basa en resolver el gradiente nulo y hessiana definida negativa.

* ¿Qué dificultad se presenta con el aumento de la dimensiones N de parámetros de la función de recompensa? ¿Qué métodos permiten solucionarlo?

El número de ecuaciones aumenta orden N para el gradiente y N^2 para la hessiana.

Los algoritmos genéticos permiten ir explorando los valores del espacio de parámetros inteligentemente.

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

PROBLEMA 3:

* ¿Cómo se resuelve analíticamente este problema?

La búsqueda de máximos de funciones multivariantes se basa en resolver el gradiente nulo y hessiana definida negativa.

* ¿Qué dificultad se presenta con el aumento de la dimensiones N de parámetros de la función de recompensa? ¿Qué métodos permiten solucionarlo?

El número de ecuaciones aumenta orden N para el gradiente y N^2 para la hessiana.

Los algoritmos genéticos permiten ir explorando los valores del espacio de parámetros inteligentemente.

* Suponga que es capaz, de encontrar los parámetros que cumplen la optimización. ¿Puede asegurar que es la política óptima (es la que maximiza la recompensa)?

Política

$$(r, p)^* = \operatorname{argmax}_{(r, p)} R = \operatorname{argmin}_{(r, p)} T$$

PROBLEMA 3:

* ¿Cómo se resuelve analíticamente este problema?

La búsqueda de máximos de funciones multivariantes se basa en resolver el gradiente nulo y hessiana definida negativa.

* ¿Qué dificultad se presenta con el aumento de la dimensiones N de parámetros de la función de recompensa? ¿Qué métodos permiten solucionarlo?

El número de ecuaciones aumenta orden N para el gradiente y N^2 para la hessiana.

Los algoritmos genéticos permiten ir explorando los valores del espacio de parámetros inteligentemente.

* Suponga que es capaz, de encontrar los parámetros que cumplen la optimización. ¿Puede asegurar que es la política óptima (es la que maximiza la recompensa)?

No, pues solo nos estamos centrando en un subconjunto del espacio de políticas.

π

POLÍTICA es una función que para la observación actual del ambiente devuelve la acción a tomar.

PROBLEMA 4: ¿Puede la Política depender de las observaciones previas? ¿En qué caso pasa?

POLÍTICA es una función que para la observación actual del ambiente devuelve la acción a tomar.

PROBLEMA 4: ¿Puede la Política depender de las observaciones previas? ¿En qué caso pasa?

Cuando una observación no tenga toda la información del sistema, existen variables ocultas. En este caso, necesitaremos hallar las variables ocultas usando información de otras observaciones.

Ejemplo: Sea un sistema caracterizado por la posición y velocidad (x,v) de un objeto; pero, cada observación solo es sobre la posición. En este caso, la variable oculta velocidad v_t podemos calcularla como $x_t - x_{t-1}$.

La función política dependerá de las observaciones actual y anterior.

POLÍTICA es una función que para la observación actual del ambiente devuelve la acción a tomar.

PROBLEMA 4: ¿Puede la Política depender de las observaciones previas? ¿En qué caso pasa?

Cuando una observación no tenga toda la información del sistema, existen variables ocultas. En este caso, necesitaremos hallar las variables ocultas usando información de otras observaciones.

Ejemplo: Sea un sistema caracterizado por la posición y velocidad (x,v) de un objeto; pero, cada observación solo es sobre la posición. En este caso, la variable oculta velocidad v_t podemos calcularla como $x_t - x_{t-1}$.

La función política dependerá de las observaciones actual y anterior.

* ¿Qué pasa cuando la observación es ruidosa?

POLÍTICA es una función que para la observación actual del ambiente devuelve la acción a tomar.

PROBLEMA 4: ¿Puede la Política depender de las observaciones previas? ¿En qué caso pasa?

Cuando una observación no tenga toda la información del sistema, existen variables ocultas. En este caso, necesitaremos hallar las variables ocultas usando información de otras observaciones.

Ejemplo: Sea un sistema caracterizado por la posición y velocidad (x,v) de un objeto; pero, cada observación solo es sobre la posición. En este caso, la variable oculta velocidad v_t podemos calcularla como $x_t - x_{t-1}$.

La función política dependerá de las observaciones actual y anterior.

* ¿Qué pasa cuando la observación es ruidosa?

Se utiliza métodos de inferencia para la determinación del estado del sistema.

POLÍTICA es una función que para la observación actual del ambiente devuelve la acción a tomar.

PROBLEMA 4: ¿Puede la Política depender de las observaciones previas? ¿En qué caso pasa?

Cuando una observación no tenga toda la información del sistema, existen variables ocultas. En este caso, necesitaremos hallar las variables ocultas usando información de otras observaciones.

Ejemplo: Sea un sistema caracterizado por la posición y velocidad (x,v) de un objeto; pero, cada observación solo es sobre la posición. En este caso, la variable oculta velocidad v_t podemos calcularla como $x_t - x_{t-1}$.

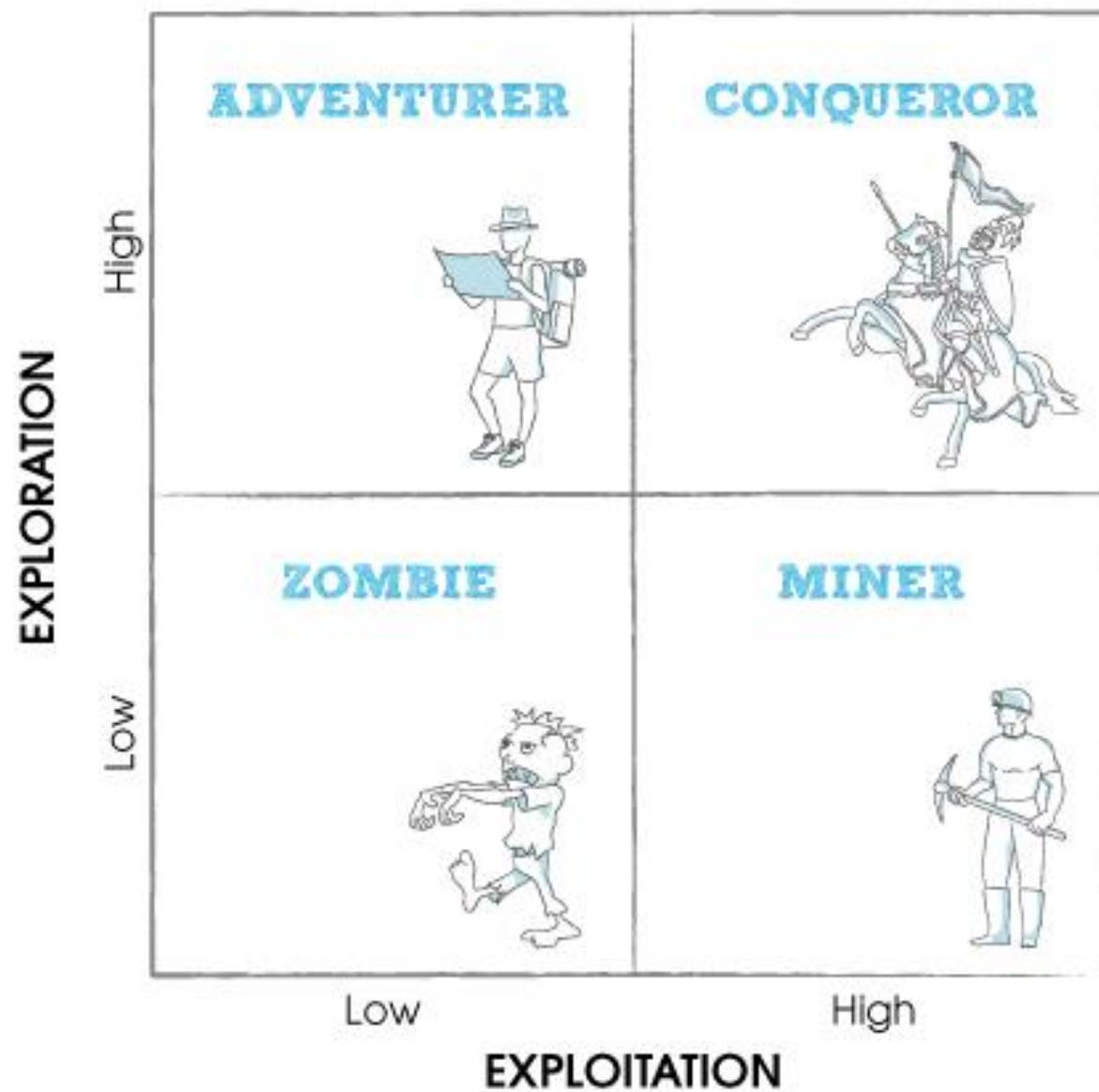
La función política dependerá de las observaciones actual y anterior.

* ¿Qué pasa cuando la observación es ruidosa?

Se utiliza métodos de inferencia para la determinación del estado del sistema.

OBSERVACIÓN = ESTADO COMPLETO DEL SISTEMA

POLÍTICA es una función que para la observación actual del ambiente devuelve la acción a tomar.



π

Procesos de decisión markoviana.

π

Procesos de decisión markoviana.

Un MDP está definido por:

π

Procesos de decisión markoviana.

Un MDP está definido por:

- Un conjunto de estados S .

π

Procesos de decisión markoviana.

Un MDP está definido por:

- Un conjunto de estados S .
- Una condición inicial $s_{t=0} \in S$.

π

Procesos de decisión markoviana.

Un MDP está definido por:

- Un conjunto de estados S .
- Una condición inicial $s_{t=0} \in S$.
- Un conjunto de acciones A .

π

Procesos de decisión markoviana.

Un MDP está definido por:

- Un conjunto de estados S .
- Una condición inicial $s_{t=0} \in S$.
- Un conjunto de acciones A .
- Una dinámica de transición $P(s'|s,a)$

Procesos de decisión markoviana.

Un MDP está definido por:

- Un conjunto de estados S .
- Una condición inicial $s_{t=0} \in S$.
- Un conjunto de acciones A .
- Una dinámica de transición $P(s'|s,a)$
- Una función de recompensa $R = R(s,a,s')$

Procesos de decisión markoviana.

Un MDP está definido por:

- Un conjunto de estados S .
- Una condición inicial $s_{t=0} \in S$.
- Un conjunto de acciones A .
- Una dinámica de transición $P(s'|s,a)$
- Una función de recompensa $R = R(s,a,s')$

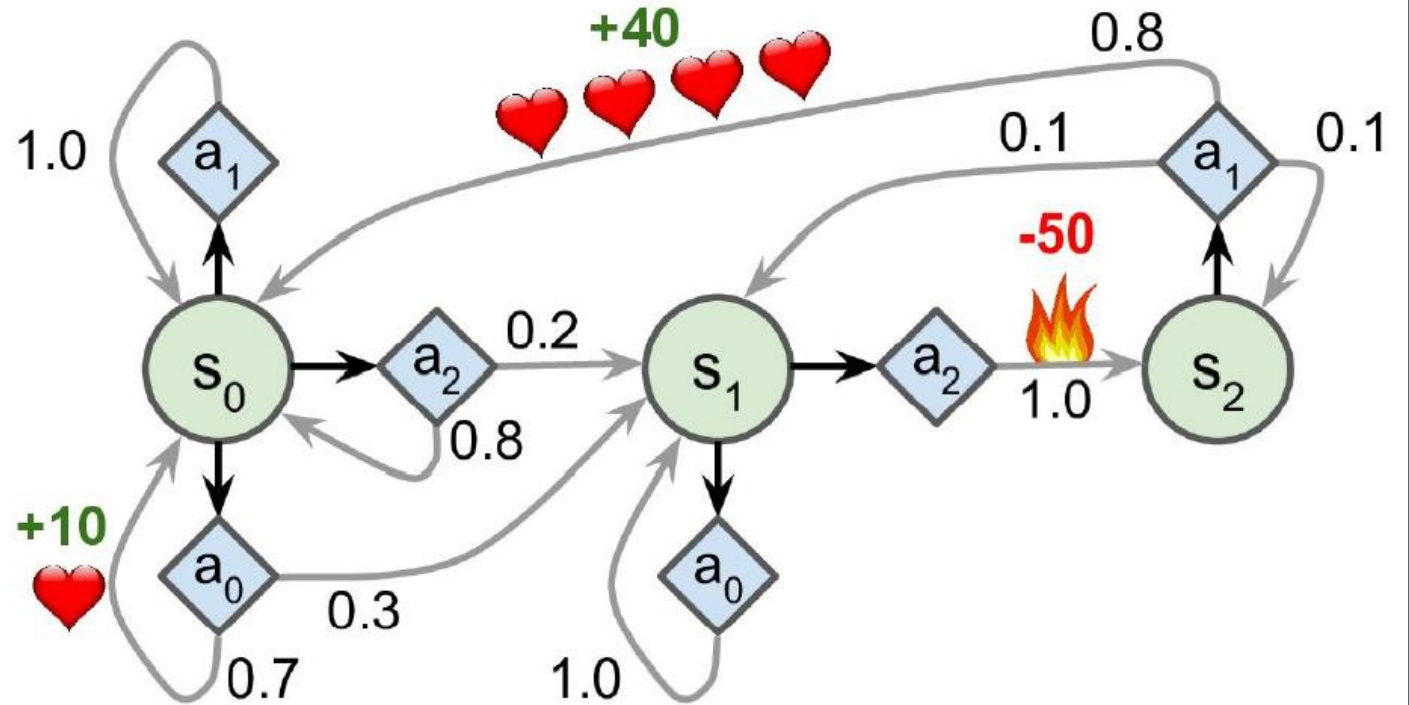
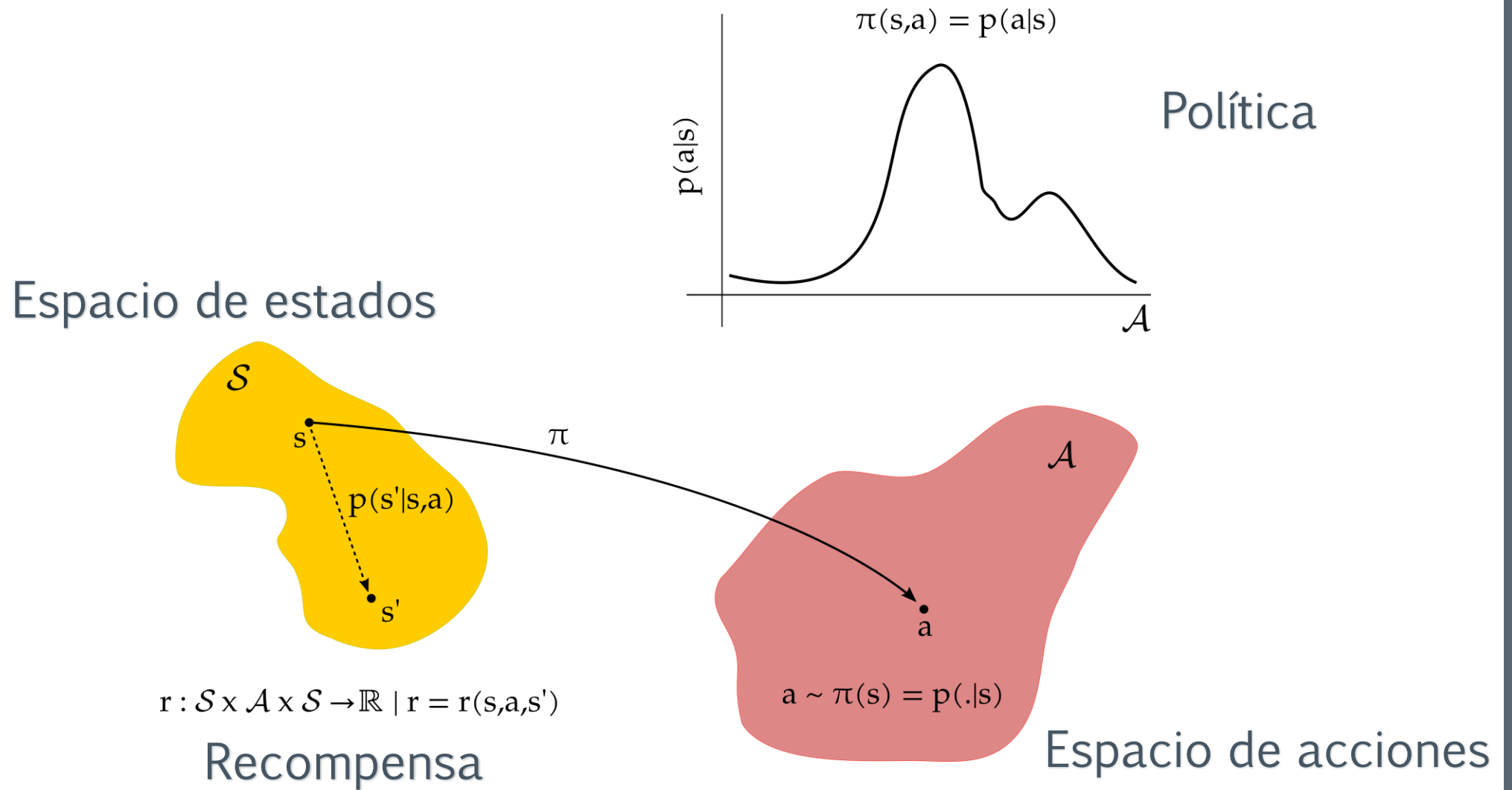


Figure 16-8. Example of a Markov decision process

π

Procesos de decisión markoviana.

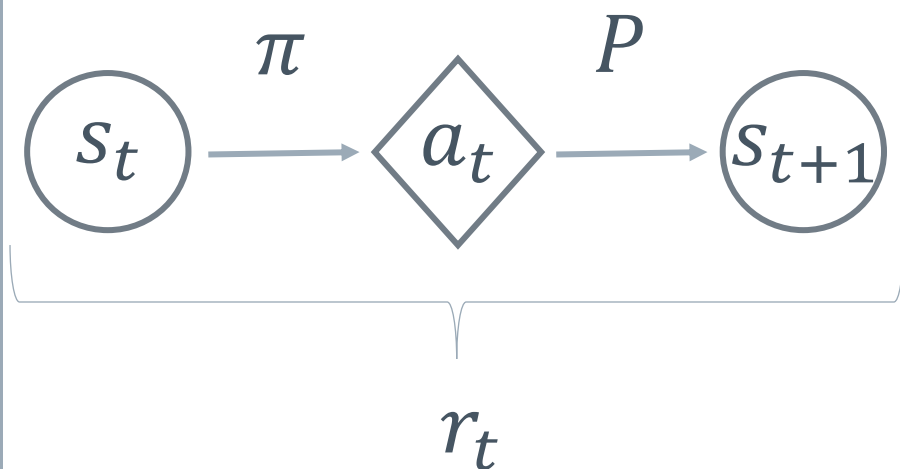


π

Procesos de decisión markoviana.

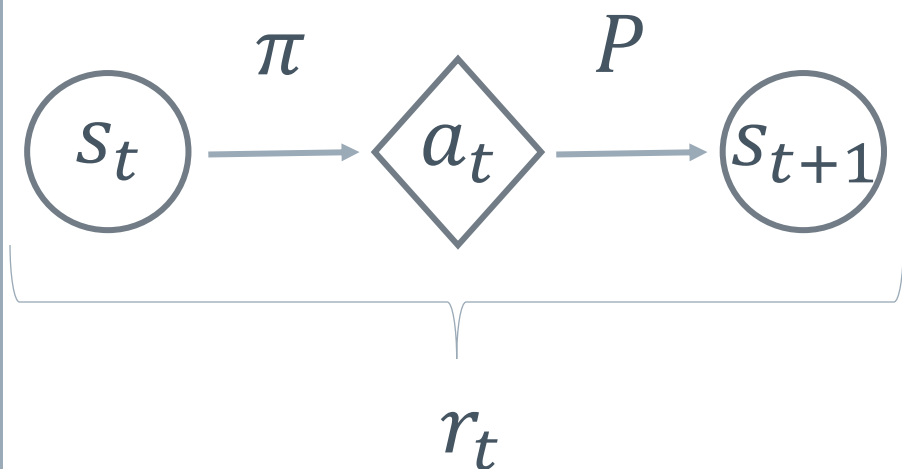
π

Procesos de decisión markoviana.



π

Procesos de decisión markoviana.

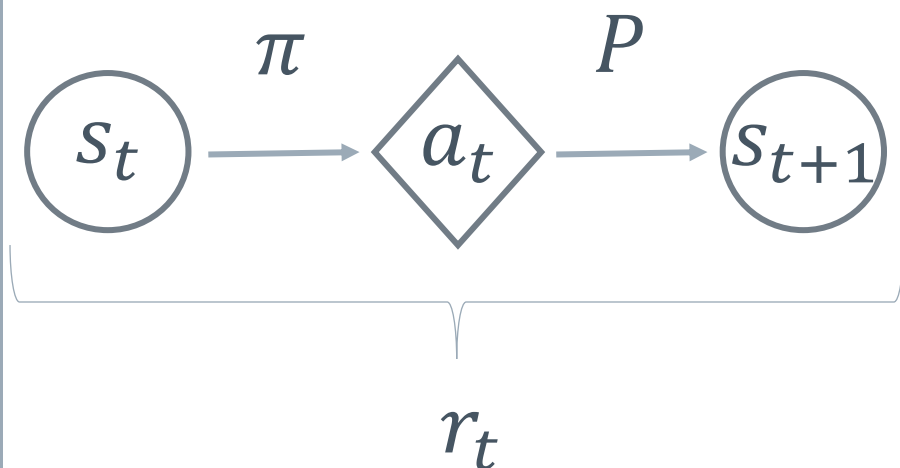


Recompensa acumulada
(variable aleatoria)

$$R_t^\pi = \sum_{n=0}^{\infty} \gamma^n r_{t+n}$$

π

Procesos de decisión markoviana.



Recompensa acumulada
(variable aleatoria)

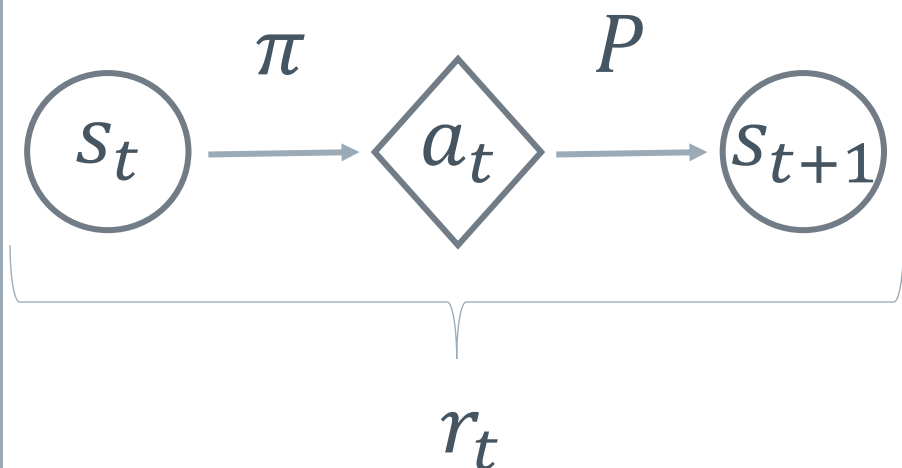
$$R_t^\pi = \sum_{n=0}^{\infty} \gamma^n r_{t+n}$$

Value function

$$V^\pi(s) = \mathbb{E}[R_t^\pi | s_t = s]$$

π

Procesos de decisión markoviana.



Recompensa acumulada
(variable aleatoria)

$$R_t^\pi = \sum_{n=0}^{\infty} \gamma^n r_{t+n}$$

Value function

$$V^\pi(s) = E[R_t^\pi | s_t = s]$$

PROBLEMA: Deducir la ecuación de Bellman (1)

$$V^\pi(s) = E[r(s, \pi(s), s')] + \gamma \sum_{s'} p(s'|s, \pi(s)) V^\pi(s')$$

π

Procesos de decisión markoviana.

Q function

$$Q^{\pi}(s, a) = E[R_t^{\pi} | s_t = s, a_t = a]$$

Procesos de decisión markoviana.

Q function

$$Q^{\pi}(s, a) = E[R_t^{\pi} | s_t = s, a_t = a]$$

PROBLEMA: Deducir la ecuación de Bellman (2)

$$Q^{\pi}(s, a) = E[r(s, a, s')] + \gamma \sum_{s'} p(s'|s, a) V^{\pi}(s')$$

Procesos de decisión markoviana.

Q function

$$Q^{\pi}(s, a) = E[R_t^{\pi} | s_t = s, a_t = a]$$

PROBLEMA: Deducir la ecuación de Bellman (2)

$$Q^{\pi}(s, a) = E[r(s, a, s')] + \gamma \sum_{s'} p(s'|s, a) V^{\pi}(s')$$

PROBLEMA: Deducir la relación

$$Q^{\pi}(s, \pi(s)) = V^{\pi}(s)$$

π

Política óptimas

Una política π es óptima si maximiza V

$$V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$$

π

Política óptimas

Una política π es óptima si maximiza V

$$V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$$

$$Q^*(s, \pi^*(s)) = V^*(s)$$

π

Política óptimas

Una política π es óptima si maximiza V

$$V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$$

$$Q^*(s, \pi^*(s)) = V^*(s)$$

Equivalentemente,

$$\max_a Q^*(s, a) = V^*(s)$$

π

Política óptimas

Una política π es óptima si maximiza V

$$V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$$

$$Q^*(s, \pi^*(s)) = V^*(s)$$

Equivalentemente,

$$\max_a Q^*(s, a) = V^*(s)$$

Entonces, la política óptima se puede calcular

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

Política óptimas

PROBLEMA: Desde (2) deducir la ecuación de Bellman (3)

$$Q^*(s, a) = E[r(s, a, s')] + \gamma \sum_{s'} p(s'|s, a) \max_{a'} Q^*(s', a')$$

Política óptimas

PROBLEMA: Desde (2) deducir la ecuación de Bellman (3)

$$Q^*(s, a) = E[r(s, a, s')] + \gamma \sum_{s'} p(s'|s, a) \max_{a'} Q^*(s', a')$$

Si resolvemos la ecuación anterior (3), la política óptima es

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

Algoritmos del RL

$$Q^*(s, a) = E[r(s, a, s')] + \gamma \sum_{s'} p(s'|s, a) \max_{a'} Q^*(s', a')$$

Programación dinámica.

(Value Iteration, Q Iteration)

Necesito conocer el MDP completo

Monte Carlo

Estimación de la dinámica de transición P y de la Recompensa acumulada.

No explota la propiedad de Markov

Temporal Difference

(Q-learning, SARSA)

Estimación de P

Explota la propiedad de Markov

Algoritmos del RL: Temporal Difference.

SARSA (On Policy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \underline{Q(S', A')} - Q(S, A)]$$

Q-Learning (Off Policy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \underline{\max_a Q(S', a)} - Q(S, A)]$$

Algoritmos del RL y Deep Learning.

$$Q: S \times A \rightarrow R$$

	a_0	a_1	a_2	...		a_M
s_0						
s_1						
s_2						
s_N						

$$Q \in R^{|S| \times |A|}$$

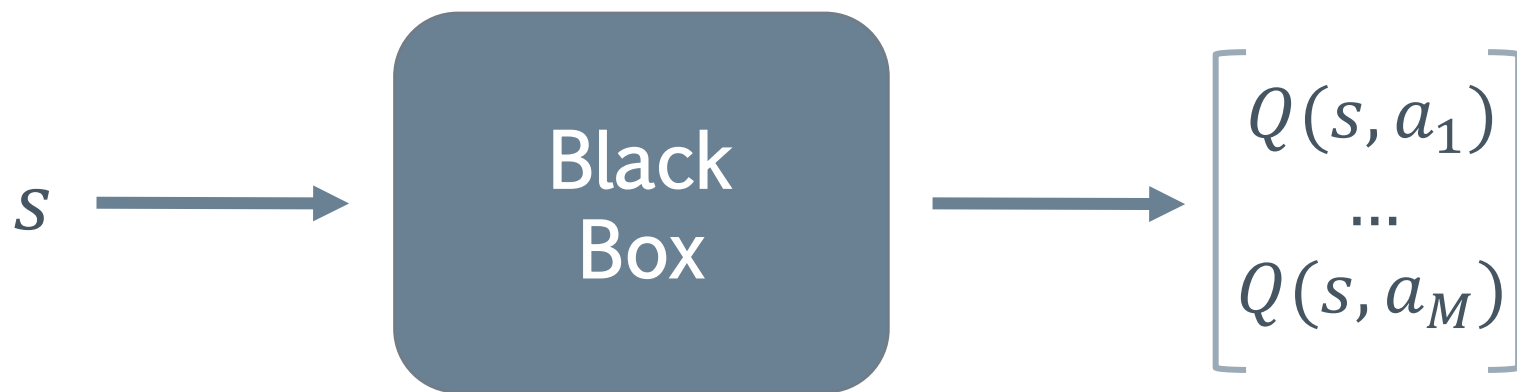
$$|A| = M$$

$$|S| = N$$

Algoritmos del RL y Deep Learning.

$$\begin{aligned} |A| &= M \\ |S| &\rightarrow \infty \end{aligned}$$

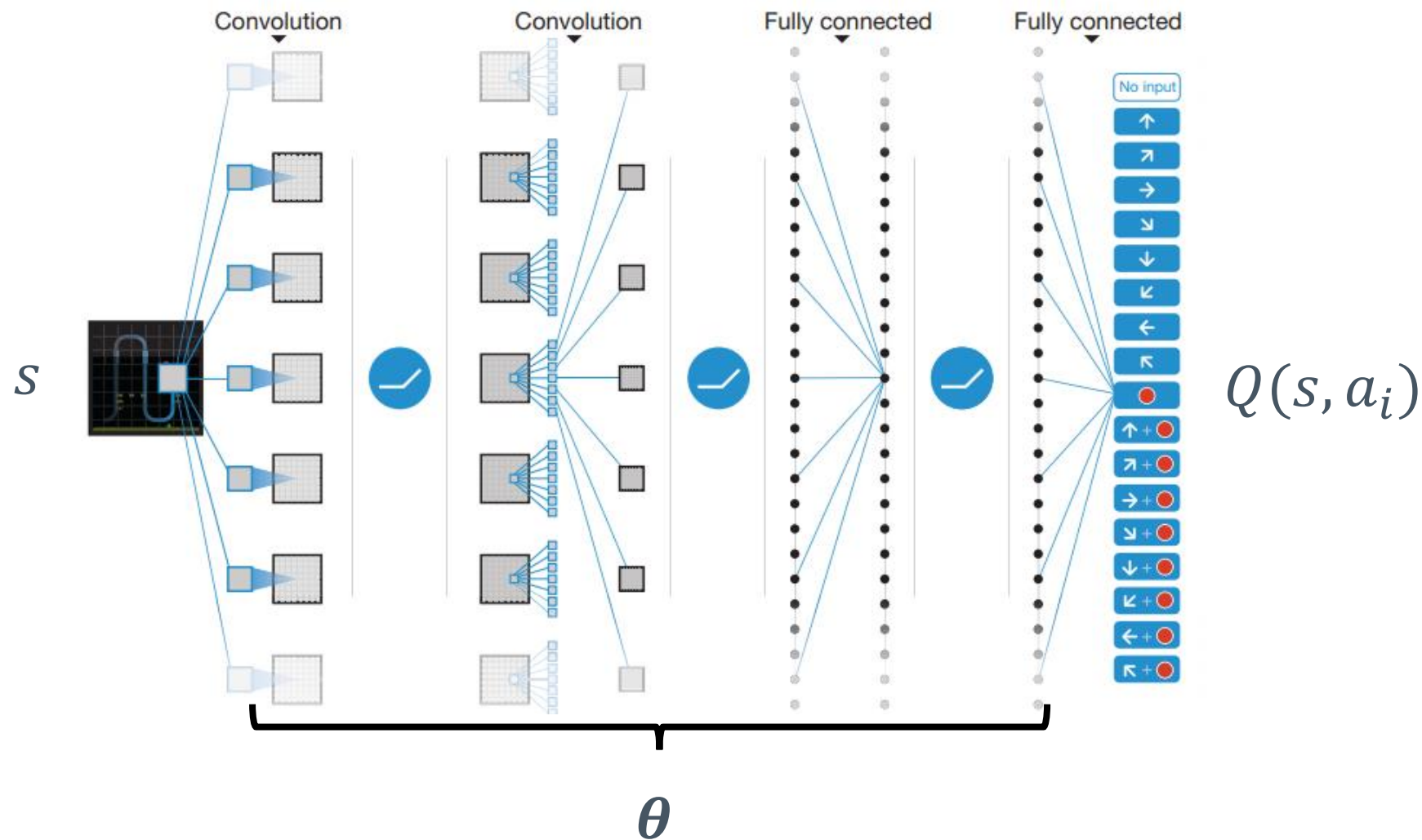
$$Q = \begin{bmatrix} Q(s, a_1) \\ \dots \\ Q(s, a_M) \end{bmatrix}$$



π

$$|A| = M$$

$$|S| \rightarrow \infty$$



$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{U}(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right]$$

Deep Q-Learning