

APRENDIZAJE REFORZADO

CLASE 7

Julían Martínez

¿QUÉ HACER CUÁNDO EL ESPACIO DE ESTADOS ES MUY GRANDE/ CONTINUO?

Reinforcement learning can be used to solve **large** problems, e.g.

- ▶ Backgammon: 10^{20} states
- ▶ Go: 10^{170} states
- ▶ Helicopter: continuous state space
- ▶ Robots: informal state space (physical universe)

How can we scale up our methods for **prediction** and **control**?

APROXIMACIÓN DE FUNCIÓN DE VALOR

Idea principal: Cambiar el aproximar la función de valor *para cada estado* (approach *tabular*), por **aproximarla globalmente**.

$$v_{\pi}(s) \approx \hat{v}(s; w)$$

EL MENÚ PARA $\hat{v}(s; w)$

Casi toda herramienta que sirve para **Aprendizaje Supervisado**:

- Árboles de decisión
- Regresión lineal
- Redes neuronales
- Fourier

FUNCIÓN DE COSTO/MÉTODO DEL GRADIENTE

$$J(\omega) := E_{\mu}[(r_{\pi}(S) - \hat{r}(S, \omega))^2]$$

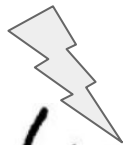
Típicamente μ es la medida para p_{π}^{π} .

$$\Delta \omega = -\frac{1}{2} \alpha \nabla_{\omega} J(\omega)$$

$$= \alpha E_{\mu}[(r_{\pi}(S) - \hat{r}(S, \omega)) \cdot \nabla_{\omega} \hat{r}(S, \omega)]$$

STOCHASTIC GRADIENT - INCREMENTAL

$$\Delta \omega = \alpha (v_{\pi}(S) - \hat{v}(S, \omega)) \cdot \nabla_{\omega} \hat{v}(S, \omega)$$



$$\omega^{k+1} = \omega^k - \alpha (v_{\pi}(S_k) - \hat{v}(S_k, \omega^k)) \cdot \nabla_{\omega} \hat{v}(S_k, \omega^k)$$

$$\Delta \omega := \alpha (G_t - \hat{v}(S_t, \omega)) \cdot \nabla_{\omega} \hat{v}(S_t, \omega)$$

$$\Delta \omega = \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, \omega) - \hat{v}(S_t, \omega)) \cdot \nabla_{\omega} \hat{v}(S_t, \omega)$$

Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated

Input: a differentiable function $\hat{v} : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameter: step size $\alpha > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop forever (for each episode):

 Generate an episode $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$ using π

 Loop for each step of episode, $t = 0, 1, \dots, T - 1$:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$

Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated

Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$

Algorithm parameter: step size $\alpha > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose $A \sim \pi(\cdot | S)$

 Take action A , observe R, S'

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})$

$S \leftarrow S'$

 until S is terminal

CASO SIMPLE E IMPORTANTE - APROXIMACIÓN LINEAL

$$\phi(s) = \begin{bmatrix} \phi_1(s) \\ \vdots \\ \phi_n(s) \end{bmatrix} \quad \hat{v}(s; w) = \phi(s)^T w = \sum_{j=1}^n \phi_j(s) w_j$$

$$\delta_t := R_{t+1} + \gamma \hat{v}(S_{t+1}; w) - \hat{v}(S_t; w)$$

$$\nabla_w \hat{v}(S_t; w) = \phi(S_t) = \phi_t$$

$$w^{t+1} = w^t + \alpha \delta_t(w^t) \phi_t$$

SOBRE LA CONVERGENCIA - MC

$$\min_{\omega} E_{\mu}[(G_t - r_{\omega}(S_t))^2]$$

$$\min_{\beta} E[(Y - \beta^T X)^2]$$

$$\beta^* = E[XX^T]^{-1} \cdot E[Y \cdot X]$$

SOBRE LA CONVERGENCIA - TD

$$w^{t+1} = w^t + \alpha (R_{t+1} + \gamma w^t \phi_{t+1} - w^t \phi_t) \cdot \phi_t$$

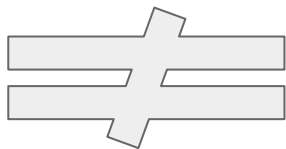
$$E[w^{t+1} | w^t = w] = w + \alpha E_w[R_{t+1} \phi_t] - \alpha w \cdot E_w[(\phi_t - \gamma \phi_{t+1})^T \phi_t]$$

$$\Rightarrow E_{\gamma}[R_{t+1} \cdot \phi(S_t)] = w E_{\gamma}[(\phi(S_t) - \gamma \phi(S_{t+1}))^T \phi(S_t)]$$

SOBRE CONVERGENCIA - TD

$$\min_w E[(R_{t+1} + \gamma v_w(S_{t+1}) - v_w(S_t))^2]$$

$$w^{t+1} = w^t + \alpha (R_{t+1} + \gamma w^t \phi_{t+1} - w^t \phi_t) \cdot (\phi_t - \gamma \phi_{t+1})$$



$$w^{t+1} = w^t + \alpha (R_{t+1} + \gamma w^t \phi_{t+1} - w^t \phi_t) \cdot \phi_t$$

COMO USAR LOS DATOS MÁS VECES (HECHARLE AGUA AL TANG)

$$\mathcal{D} = \{(s_1, \hat{v}_1^\pi), \dots, (s_T, \hat{v}_T^\pi)\} \quad \text{Experience Replay (BATCH)}$$

$$\hat{v}_i^\pi \approx v_\pi(s_i)$$

$$\arg\min_w \sum_{i=1}^T (\hat{v}_i^\pi - \hat{v}_w(s_i))^2$$

$$\arg\min_w \mathbb{E}_{p_T} [(\hat{v}^\pi(s) - \hat{v}_w(s))^2]$$

EJEMPLO CONCRETO DE PROBLEMAS CON LA NOTACIÓN

Stochastic Gradient Descent with Experience Replay

Given experience consisting of $\langle \text{state}, \text{value} \rangle$ pairs

$$\mathcal{D} = \{\langle S_1, \hat{v}_1^\pi \rangle, \langle S_2, \hat{v}_2^\pi \rangle, \dots, \langle S_T, \hat{v}_T^\pi \rangle\}$$

Repeat:

1. Sample state, value from experience

$$\langle s, \hat{v}^\pi \rangle \sim \mathcal{D}$$

2. Apply stochastic gradient descent update

$$\Delta\theta = \alpha(\hat{v}^\pi - v_\theta(s))\nabla_\theta v_\theta(s)$$

Converges to least squares solution

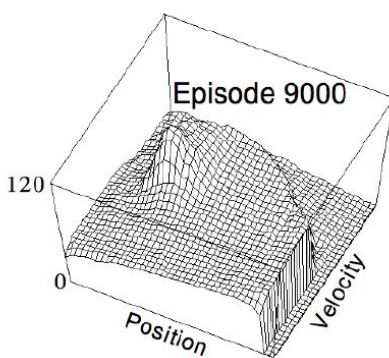
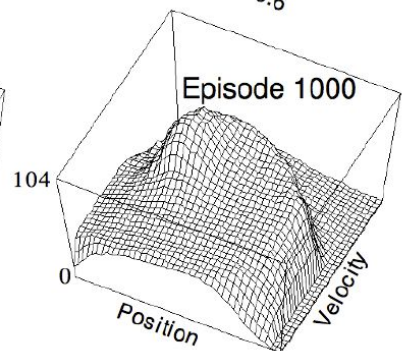
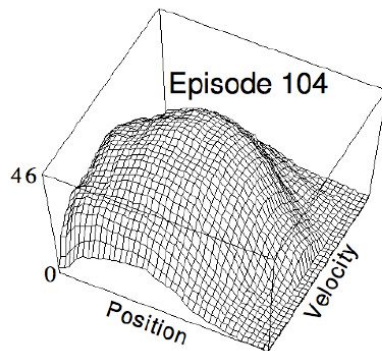
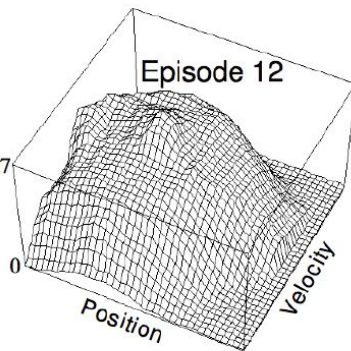
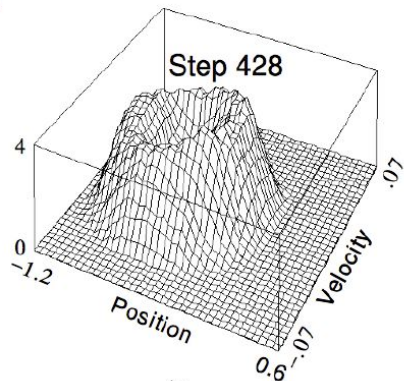
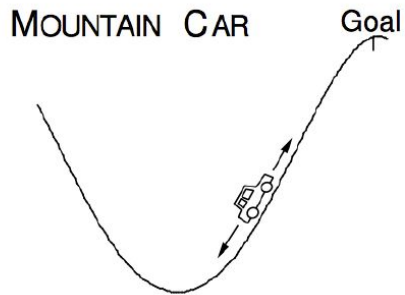
$$\theta^\pi = \underset{\theta}{\operatorname{argmin}} \operatorname{LS}(\theta) = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{\mathcal{D}} [(\hat{v}_i^\pi - v_\theta(S_i))^2]$$

LINEAR LEAST SQUARES PREDICTION

$$\mathbf{w}_B^* = \left(\sum_{t=1}^T \phi(s_t) \cdot \phi(s_t)^T \right)^{-1} \sum_{t=1}^T \phi(s_t) \hat{v}_t^\pi$$

Ver algoritmo Sherman-Morrison para resolver esto con menor complejidad.

SOBRE LA "GRANULARIDAD DEL APROXIMANTE"



EL LADO OSCURO...



- Sobre la convergencia de temporal difference
<https://papers.nips.cc/paper/3809-convergent-temporal-difference-learning-with-arbitrary-smooth-function-approximation.pdf>
- Detalles y referencias sobre gradient temporal difference y convergencia
<https://towardsdatascience.com/reinforcement-learning-an-introduction-to-gradient-temporal-difference-learning-algorithms-4a72ce5ab31e>

