

# Aprendizado Por Reforço

Introdução

Framework

Histórico e Estados

Outros Componentes

Markov Reward Process

Função de Valor

Markov Decision Process

Política

Função de Valor

Programa de Backup

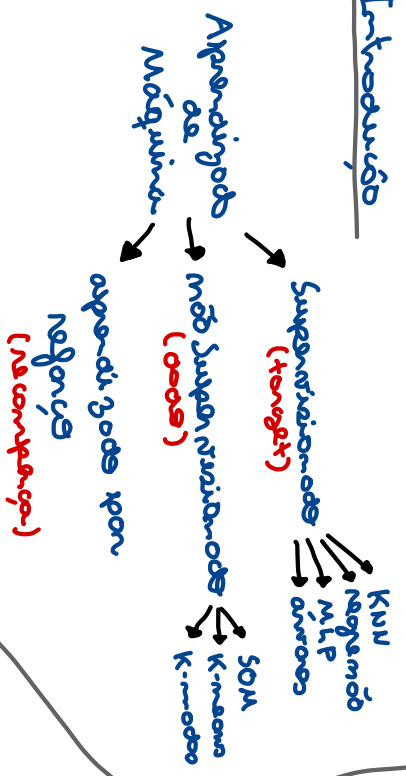
Função de Valor Ótima

Estratégias de Aprendizado

Programação Dinâmica

Q-learning

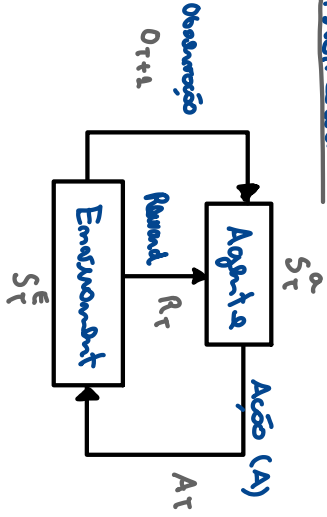
# Introdução



Exemplos:

- \* Self drive car
- \* Chat bot
- \* Jogo

## Framework



Goal: Maximizar a recompensa acumulada

Exemplos:

- \* Atari
  - + Reward na captura o score
  - Reward na diminuição o score
- \* Finanças
  - + Reward de ganhar dinheiro
  - Reward de perder dinheiro
- \* Power Station
  - + Reward de produzir energia de forma segura
  - Reward de acabar a produção segura de reactor

## Histórias e Estados

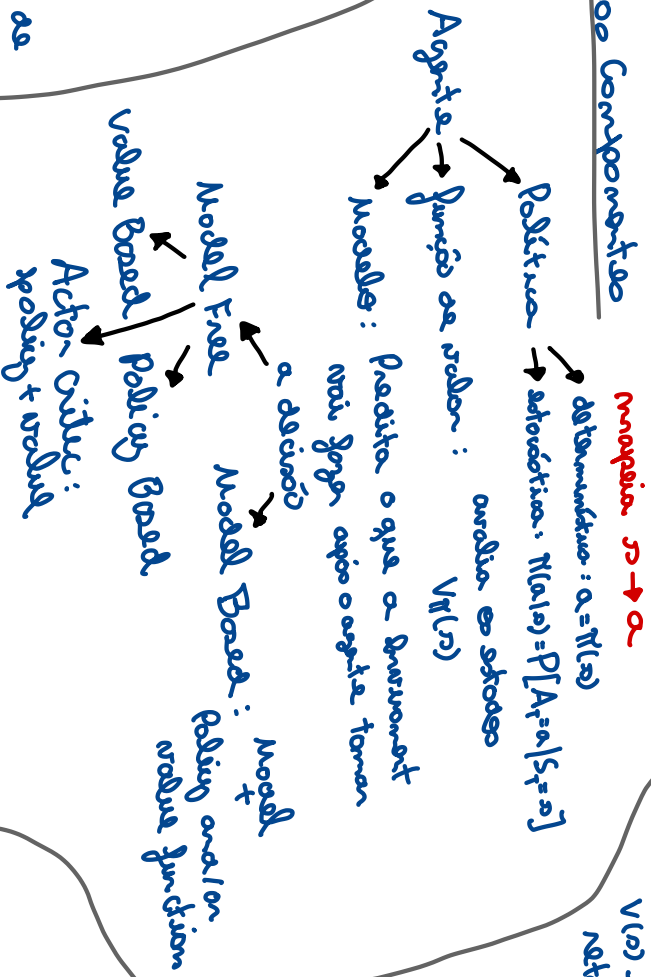
$H_t = [(O_t, A_t, R_t), \dots, (O_t, A_t, R_t)]$   
 $S_t = f(H_t)$   
 Quando usa toda a história?

- \* Mais pontos do problema em RL são modelados como codição de Markov

Definição:  $S_t$  e Markov na  $P[S_{t+1} | S_t, S_{t-1}, \dots, S_1] = P[S_{t+1} | S_t]$

## Não usam toda a história e recompensa

### Outros Componentes



## Markov Reward Process

- É uma tupla  $\langle S, P, R, \gamma \rangle$
- \* S: States
  - \* P: Matriz de transição
  - \* R: Função de recompensa,  $R_s = E[R_{t+1} | S_t = s]$
  - \*  $\gamma$ : fator de desconto,  $\gamma \in [0, 1]$

### Retorno:

total descontada reward

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

### Função de Valor

$V(s)$  é o valor esperado de retornos iniciando no estado  $s$

$$V(s) = E[G_t | S_t = s]$$

## Função de Valor (cont.)

Seja  $G_\tau = \sum_{k=0}^{\infty} \gamma^k R_{\tau+k+1}$  e  $V(s) = E[G_\tau | S_\tau = s]$ ,

a função de valor pode ser decomposta

em duas partes:

- \* Recompensa Imediata
- \* Valores com desconto das recompensas futuras

Logo:

$$V(s) = E[G_\tau | S_\tau = s]$$

$$E[R_{\tau+1} + \gamma R_{\tau+2} + \gamma^2 R_{\tau+3} + \dots | S_\tau = s]$$

$$E[R_{\tau+1} + \gamma (R_{\tau+2} + \gamma R_{\tau+3} + \dots) | S_\tau = s]$$

$$E[R_{\tau+1} + \gamma V(s) | S_\tau = s]$$

Imediatamente  
decomposta  
em  
recompensa  
e valor futuro

Chamamos isso de Bellman Equation:

$$V(s) = R_s + \gamma \sum_{s'} P_{ss'} V(s')$$

Para resolver basta:

$$V = (I - \gamma P)^{-1} R$$

\* Complexidade  $O(n^3)$

Monta matriz

\* Método Iterativo

Programação dinâmica

Temporal-Difference Learning

## Markov Decision Process

\* MDP é uma MRP com decisão (actions)

\* MDP é uma tupla  $\langle S, A, P, R, \gamma \rangle$

\* S: conjunto de estados

A: conjunto de ações

P: matriz de transição

R: função de recompensa  $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$

$\gamma$ : fator de desconto,  $\gamma \in [0, 1]$

## Política

$$\pi(a|s) = P[A_t = a | S_t = s]$$

Função de valor:

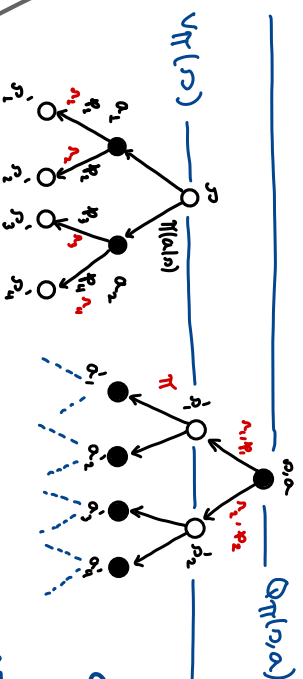
Seja a função de valor: seguindo uma política  $\pi$ ,  $s$

$$V_\pi(s) = E[G_\tau | S_\tau = s]$$

Action value function: seguindo uma política  $\pi$ ,  $s, a$

$$q_\pi(s, a) = E[G_\tau | S_\tau = s, A_\tau = a]$$

Programo de backup: Ajuda x0 manter as equações



$$q_\pi(s, a) = \sum_{s'} P(s'|s, a) [r + \gamma V_\pi(s')]$$

$$V_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$$

## Função de Valor Ótima

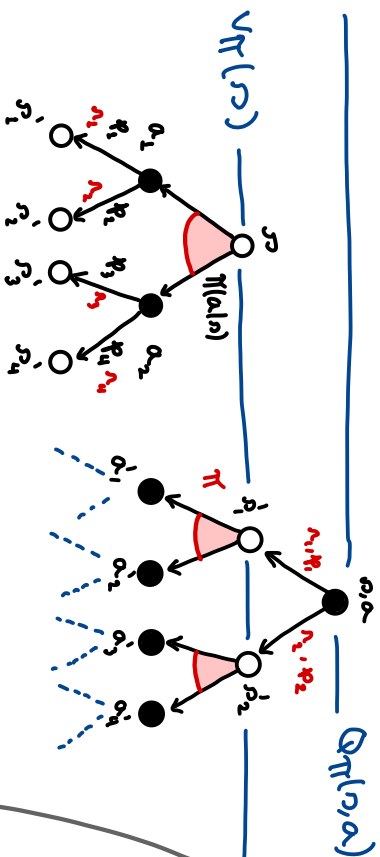
Existe uma política que maximiza a função de valor:

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

A melhor política pode ser encontrada a partir de  $Q(s, a)$

$$\pi_* = \begin{cases} 1 & \text{se } a = \arg\max (Q_*(s, a)) \forall a \in A \\ 0 & \text{caso contrário} \end{cases}$$



$$V_*(s) = \max \{Q_*(s, a)\}$$

## Como encontrar a melhor Política

Progamção dinâmica → Policy Iteration  
 métodos → Value Iteration

RL Algp. → TD-Learning  
 → Q-Learning

### Programação Linear

- \* Divide um problema em sub problemas
- \* É recursivo
- \* Preciso conhecer o Environment (Model board)
- \* Usa Bellman eq. properties

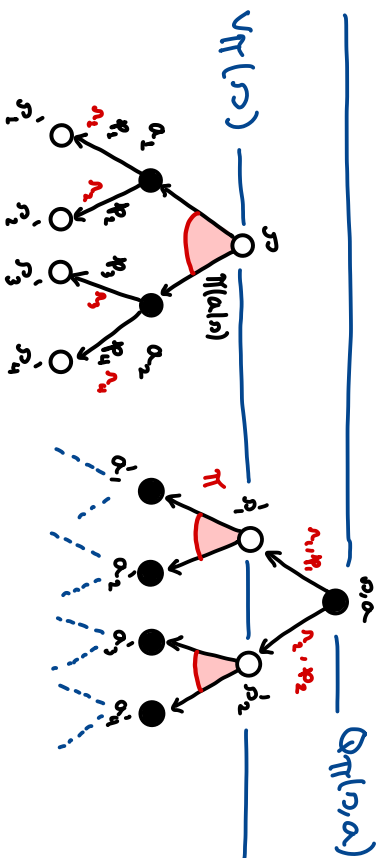
### Q-Learning

- \* Exploration & exploitation
- \* Uno  $\Theta Q(s, a)$  para encontrar a melhor ação
- \* Model Free, não preciso conhecer o environment
- \* mais simples em RL.

# Programação

# Dinâmica

## Grid World



while  $\Delta > \epsilon$ :  
 calculate  $V$   
 $\Delta = |V(s) - V'(s)|$

## Policy Iteration

$$V(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) [R + \gamma V(s')]$$

## Value Iteration

$$V(s) = \max_a \left\{ \sum_{s'} p(s'|s,a) [R + \gamma V(s')] \right\}$$

Policy  $\times$  value  
 iteration  
 converge to optimal  $\pi$   
 more interactions

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16