

Métodos Monte Carlo

Na aula passada vimos métodos que requerem conhecimento da dinâmica do processo de decisão de Markov para resolvê-lo. Nessa aula iremos estudar métodos que aprendem com a experiência do agente no ambiente. Ou seja, não é necessário ter conhecimento da função de transição $p(s', r|s, a)$.

Mesmo se for possível obter a função $p(s', r|s, a)$, muitas vezes o seu cálculo não é trivial. Por exemplo, qual a probabilidade de aparecer um Rei de Ouros na próxima rodada de uma mesa de pôquer? Um programa de computador pode realizar o cálculo, mas a sua implementação poder ser trabalhosa e sujeita a erros.

Previsão com Monte Carlo

Dada uma política π , podemos aproximar os valores de V através de simulações de episódios com π :

$$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T.$$

O valor de $V(s_t)$ é a média dos valores de recompensa acumulada G_t obtida após estado s ser visitado.

Existem duas possíveis implementações dessa ideia:

- Primeira visita: a média é calculada para os valores G_t da primeira visita a um estado s em um episódio. Iremos focar nessa implementação nessa aula.
- Todas visitas: a média é calculada para todos os valores G_t de todas as visitas a um estado em um episódio. Essa ideia é útil quando utilizando **traços de elegibilidade**.

O algoritmo abaixo mostra o pseudocódigo de um método MC para previsão utilizando a estratégia de “primeira visita”. O código abaixo assume que o agente interage infinitamente com o ambiente, obtendo episódios que são utilizados para estimar $V(s)$.

```
1 def Primeira-Visita MC Previsão( $\pi$ ):
2   inicialize  $V(s)$  de forma arbitrária.
3   ( $V(s) = 0$  para estados terminais)
4   inicializa Retorno[ $s$ ] = [] para todos  $s$ 
5   while True:
6     gera um episódio com  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
7      $G = 0$ 
8     for  $t$  in range( $T-1, 0$ ):
9        $G = \gamma G + R_{t+1}$ 
10      if  $S_t$  not in  $\{S_0, S_1, \dots, S_{t-1}\}$ :
11        Retorno[ $S_t$ ].append( $G$ )
12         $V(S_t) = \text{média}(\text{Retorno}[S_t])$ 
```

Implementação Incremental

A implementação acima não é eficiente em termos de tempo e memória, pois requer o armazenamento de todos os valores de G para cada estado. Além disso, o algoritmo calcula a média desses valores de G toda vez que um novo valor é adicionado à lista. Uma forma de solucionar esse problema é fazendo uma implementação incremental das médias.

$$\begin{aligned} V_{k+1} &= \frac{1}{k} \sum_{i=1}^k G_i \\ &= \frac{1}{k} \left(G_k + \sum_{i=1}^{k-1} G_i \right) \\ &= \frac{1}{k} \left(G_k + (k-1) \frac{1}{(k-1)} \sum_{i=1}^{k-1} G_i \right) \\ &= \frac{1}{k} \left(G_k + (k-1) V_k \right) \\ &= \frac{1}{k} (G_k + k V_k - V_k) \\ &= V_k + \frac{1}{k} (G_k - V_k) \end{aligned}$$

Assim, podemos utilizar a última equação acima no lugar das linhas 11 e 12 do pseudocódigo.

Monte Carlo para Estimar Valores de Ações

Se um modelo do ambiente não está disponível, então estimar os valores de pares estado-ação é muito mais útil que estimar os valores de estados. Nos algoritmos de programação dinâmica, como o Iteração de Valor, uma vez obtida uma estimativa dos valores de V , uma política pode ser extraída através do processo de *lookahead*. Na ausência de um modelo, o *lookahead* não é possível.

O procedimento para estimar os valores $q_\pi(s, a)$ é o mesmo discutido acima. A diferença é que as médias são calculadas para valores G_t obtidos após visitar estado s e tomar ação a .

O problema é que, dependendo de π , alguns pares estado-ação nunca serão visitados e seus valores não serão aproximados corretamente. Uma forma de lidar com esse problema é através dos **inícios exploratórios**: todo par estado-ação pode ser escolhido como início de um episódio com probabilidade maior que zero.

Monte Carlo para Controle

Agora que já sabemos como utilizar o método Monte Carlo para estimar os valores de $q_\pi(s, a)$, podemos utilizar o método de iteração de política para encontrar uma política ótima π_* . A ideia é alternar avaliação de uma política inicialmente arbitrária e a melhoria gulosa da mesma. A melhoria gulosa com valores de q são feitos da seguinte forma:

$$\pi(s) = \arg \max_a q(s, a),$$

o que não requer o uso de um modelo do problema.

Intuitivamente, o algoritmo acima não pode convergir para uma política sub-ótima. Se ele convergisse, então o valor de V iria convergir para o valor correto daquela política, que faria a política mudar. No entanto, ainda não existe nenhuma prova formal de que o algoritmo converge para uma política ótima.

```

1 def Controle MC com Início Exploratório( $\pi$ ):
2   inicialize  $\pi(s)$  de forma arbitrária
3   inicialize  $Q(s, a)$  de forma arbitrária.
4   inicialize  $N(s, a) = 0$ 
5   while True:
6     escolha  $S_0$  e  $A_0$  aleatoriamente de forma que todos os pares tenham
        probabilidade maior que zero.
7     gera um episódio a partir de  $S_0, A_0$  seguindo  $\pi$ :
         $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
8      $G = 0$ 
9     for  $t$  in range(T-1, 0):
10       $G = \gamma G + R_{t+1}$ 
11      if  $S_t$  not in  $\{S_0, S_1, \dots, S_{t-1}\}$ :
12         $N(S_t, A_t) = N(S_t, A_t) + 1$ 
13         $Q(S_t, A_t) = Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(G - Q(S_t, A_t))$ 
14         $\pi(S_t) = \arg \max_a Q(S_t, a)$ 

```

Monte Carlo para Controle sem Inícios Exploratórios

Os inícios exploratórios nos permitem avaliar os valores de pares estado-ação. No entanto, é pouco provável que os inícios exploratórios podem ser implementados no mundo real.

Uma forma de dispensar o uso dos inícios exploratórios é derivando uma política que assinala uma probabilidade, mesmo que pequena, de executar todas as ações. Uma dessas políticas é a ϵ -gulosa, que em um estado s ,

- executa a melhor ação a com probabilidade $1 - \epsilon + \frac{\epsilon}{|A(s)|}$
- executa uma ação aleatória com probabilidade $\frac{\epsilon}{|A(s)|}$

O algoritmo acima converge para a melhor política ϵ -gulosa, que não necessariamente é ótima, mas é próxima de ótima.

Previsão com Monte Carlo Off-Policy

Uma outra alternativa para remover o requisito de inícios exploratórios é aproximar os valores de uma **política alvo** π agindo como uma **política de comportamento** b . Assim, b pode ser uma política ϵ -gulosa, que garante que todos os estados e ações são visitados, e π pode ser uma política determinística.

O requisito é que se π executa uma ação com probabilidade maior que zero, b também tem que executar a mesma ação com probabilidade maior que zero. Do contrário não seria possível aprender os valores de v_π através de episódios gerados por b , pois tais ações nunca seriam experimentadas.

Amostragem por Importância

Em estatística, **amostragem por importância** nos permite amostrar valores de uma distribuição para aproximar o valor de uma outra distribuição. A ideia é amostrar valores de $v_c(s)$ para estimar $v_\pi(s)$.

```

1 def Controle MC On-Policy ( $\epsilon$ ):
2     inicialize  $\pi(s)$  com uma política  $\epsilon$ -gulosa arbitrária
3     inicialize  $Q(s, a)$  de forma arbitrária.
4     inicialize  $N(s, a) = 0$ 
5     while True:
6         gera um episódio com  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
7          $G = 0$ 
8         for  $t$  in range(T-1, 0):
9              $G = \gamma G + R_{t+1}$ 
10            if  $t$  not in  $\{S_0, S_1, \dots, S_{t-1}\}$ :
11                 $N(S_t, A_t) = N(S_t, A_t) + 1$ 
12                 $Q(S_t, A_t) = Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(G - Q(S_t, A_t))$ 
13                 $A^* = \arg \max_a Q(S_t, a)$ 
14                for  $a$  in  $A(S_t)$ :
15                    if  $a == A^*$ :
16                         $\pi(a|S_t) = 1 - \epsilon + \epsilon/|A(S_t)|$ 
17                    else:
18                         $\pi(a|S_t) = \epsilon/|A(S_t)|$ 

```

De uma forma geral temos:

$$\begin{aligned}
 \mathbb{E}_{X \sim P}[f(X)] &= \sum P(X)f(X) \\
 &= \sum Q(x) \frac{P(X)}{Q(X)} f(X) \\
 &= \mathbb{E}_{X \sim Q} \left[\frac{P(X)}{Q(X)} f(X) \right]
 \end{aligned}$$

Para aplicarmos a amostragem por importância no problema de controle precisamos calcular o valor de probabilidade de observarmos um episódio dada a política π e dada a política b .

$$\begin{aligned}
 Pr\{A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi\} &= \pi(A_t | S_t) p(S_{t+1} | S_t, A_t) \pi(A_{t+1} | S_{t+1}) \dots p(S_T | S_{T-1}, A_{T-1}) \\
 &= \prod_{k=1}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)
 \end{aligned}$$

A fração de amostragem por importância é dada por:

$$\rho_{t:T-1} = \frac{\prod_{k=1}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)}{\prod_{k=1}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)} = \frac{\prod_{k=1}^{T-1} \pi(A_k | S_k)}{\prod_{k=1}^{T-1} b(A_k | S_k)}$$

Embora o episódio dependa da dinâmica do ambiente, como o ambiente é idêntico, os valores de p se cancelam e a fração ρ depende apenas de π e b .

Controle com Monte Carlo Off-Policy

O algoritmo de controle de Monte Carlo Off-Policy abaixo aprende sobre uma política de comportamento b enquanto deriva uma política ótima π_* . A política b tem que ser do estilo ϵ -gulosa, para garantir que todos os estados e ações são amostrados um número infinito de vezes, para garantir convergência.

```

1 def MC Previsão Off-Policy( $\pi$ ):
2     inicialize  $Q(s,a)$  de forma arbitrária.
3     while True:
4          $b$  = política  $\epsilon$ -gulosa arbitrária
5         gera um episódio utilizando  $b$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
6          $G = 0$ 
7          $W = 1$ 
8         for  $t$  in range(T-1, 0):
9              $G = \gamma G + R_{t+1}$ 
10            if  $S_t$  not in  $\{S_0, S_1, \dots, S_{t-1}\}$ :
11                 $N(S_t, A_t) = N(S_t, A_t) + 1$ 
12                 $Q(S_t, A_t) = Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(WG - Q(S_t, A_t))$ 
13                 $W = W \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$ 

```

```

1 def MC Controle Off-Policy( $\gamma$ ):
2     inicialize  $Q(s,a)$  de forma arbitrária.
3     inicialize  $\pi$  de forma arbitrária
4      $\pi(S_t) = \arg \max_a Q(S_t, a)$ 
5     while True:
6          $b$  = política  $\epsilon$ -gulosa arbitrária
7         gera um episódio utilizando  $b$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
8          $G = 0$ 
9          $W = 1$ 
10        for  $t$  in range(T-1, 0):
11             $G = \gamma G + R_{t+1}$ 
12            if  $S_t$  not in  $\{S_0, S_1, \dots, S_{t-1}\}$ :
13                 $N(S_t, A_t) = N(S_t, A_t) + 1$ 
14                 $Q(S_t, A_t) = Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(WG - V(S_t))$ 
15                 $W = W \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$ 
16                 $\pi(S_t) = \arg \max_a Q(S_t, a)$ 

```

O método de Monte Carlo com amostragem por importância é não enviesado, no sentido estatístico de que ele converge para o valor exato de q_π da política π sendo avaliada. No entanto, a variância da abordagem pode ser muito grande, o que faz com que seja necessário um número muito grande de episódios para convergir para o valor de q_π . Veremos nas próximas aulas um outro método que é enviesado mas possui menor variância.

Vantagens dos métodos Monte Carlo com relação aos métodos de programação dinâmica.

1. Não é necessário um modelo, o agente aprende através de interações diretas com o ambiente.
2. O agente aprende através de simulações, as equações de transição do problema não são necessárias.
3. É possível focar em estados que realmente interessam o agente.