

eman ta zabal zazu



Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

Técnicas de Inteligencia Artificial

Ingeniería Informática de Gestión y Sistemas de Información

Practica 4

Reinforcement Learning

Autor(es):

Xabier Gabiña

Diego Montoya

3 de diciembre de 2024

Índice general

1. Introducción	3
2. Ejercicios	4
2.1. Iteración de valores	4
2.2. Análisis de cruce de puentes	6
2.3. Q-Learning	7
3. Resultados	9

Índice de Códigos

2.1. Iteración de valores	4
2.2. Análisis de cruce de puentes	6
2.3. Q-Learning	7

1. Introducción

2. Ejercicios

2.1. Iteración de valores

Descripción

El primer ejercicio se trata de implementar un agente que realice la iteración de valores. Un agente de Iteración de Valores toma un proceso de decisión de Markov en la inicialización y ejecuta la iteración de valores durante un número dado de iteraciones utilizando el factor de descuento proporcionado.

Implementación

```
1 class ValueIterationAgent(ValueEstimationAgent):
2     """
3         * Please read learningAgents.py before reading this.*
4
5         A ValueIterationAgent takes a Markov decision process
6         (see mdp.py) on initialization and runs value iteration
7         for a given number of iterations using the supplied
8         discount factor.
9     """
10    def __init__(self, mdp, discount = 0.9, iterations = 100):
11        """
12        Your value iteration agent should take an mdp on
13        construction, run the indicated number of iterations
14        and then act according to the resulting policy.
15
16        Some useful mdp methods you will use:
17            mdp.getStates()
18            mdp.getPossibleActions(state)
19            mdp.getTransitionStatesAndProbs(state, action)
20            mdp.getReward(state, action, nextState)
21            mdp.isTerminal(state)
22        """
23        self.mdp = mdp
24        self.discount = discount
25        self.iterations = iterations
26        self.values = util.Counter() # A Counter is a dict with default 0
27        self.runValueIteration()
28
29    def runValueIteration(self):
30        # Write value iteration code here
31        for _ in range(self.iterations):
32            new_values = self.values.copy()
33            for state in self.mdp.getStates():
34                if not self.mdp.isTerminal(state):
35                    new_values[state] = max(self.computeQValueFromValues(state, action) for action
36                    in self.mdp.getPossibleActions(state))
37            self.values = new_values
38
39    def getValue(self, state):
40        """
41        Return the value of the state (computed in __init__).
42        """
43        return self.values[state]
44
45
46    def computeQValueFromValues(self, state, action):
47        """
48        computeQValueFromValues(state, action) devuelve el valor Q del par (estado, acción)
49        dado por la función de valor dada por self.values. Nota: Recuerda que para calcular
50        el valor de un estado calcularemos los q_values (estado, acción), es decir los valores
51        de las acciones posibles para quedarnos con el mayor (max de entre los q_values).
52        """
53        q_value = 0
54        for next_state, prob in self.mdp.getTransitionStatesAndProbs(state, action):
```

```

55         q_value += prob * (self.mdp.getReward(state, action, next_state) + self.discount *
self.values[next_state])
56         return q_value
57
58
59     def computeActionFromValues(self, state):
60         """
61         computeActionFromValues(state) calcula la mejor acción de acuerdo con la función de
62         valor dada por self.values. Es decir, de acuerdo al estado siguiente al que nos lleva
63         cada acción y que mayor valor tiene. Este método llamará al siguiente.
64         """
65         if self.mdp.isTerminal(state):
66             return None
67
68         return max(self.mdp.getPossibleActions(state), key=lambda x: self.computeQValueFromValues(
state, x))
69
70
71     def getPolicy(self, state):
72         return self.computeActionFromValues(state)
73
74     def getAction(self, state):
75         "Returns the policy at the state (no exploration)."
```

```

76         return self.computeActionFromValues(state)
77
78     def getQValue(self, state, action):
79         return self.computeQValueFromValues(state, action)
```

Código 2.1: Iteración de valores

Comentarios

2.2. Análisis de cruce de puentes

Descripción

Implementación

```
1 def question2():
2     answerDiscount = 0.9
3     answerNoise = 0.01
4     return answerDiscount, answerNoise
5
6 if __name__ == '__main__':
7     print('Answers to analysis questions:')
8     import analysis
9     for q in [q for q in dir(analysis) if q.startswith('question')]:
10         response = getattr(analysis, q)()
11         print(' Question %s:\t%s' % (q, str(response)))
```

Código 2.2: Análisis de cruce de puentes

Comentarios

La solución más sencilla para este ejercicio es la de bajar el ruido por debajo de 0.01. Esto hace que las probabilidades de que el agente cometa un error y se caiga a un acantilado de puntos negativos sean muy bajas. No haría falta tocar el valor gamma ya que al tener un factor alto, lo que conseguimos, es que las recompensas futuras tengan un peso mayor lo que conviene para no acabar llegando a la casilla +1 en vez de a la +10.

2.3. Q-Learning

Descripción

En este ejercicio se trata de implementar un agente que realice el aprendizaje por refuerzo mediante Q-Learning. Q-Learning es un algoritmo de aprendizaje por refuerzo que aprende una política óptima sin conocer el modelo del entorno. Se basa en ir probando acciones, en un principio de forma aleatoria, y actualizando los valores de la función Q en función de las recompensas recibidas. A medida que se va actualizando la función Q, el agente va aprendiendo la mejor acción a tomar en cada estado y va disminuyendo el número de acciones aleatorias.

Implementación

```
1 class QLearningAgent(ReinforcementAgent):
2     """
3     Q-Learning Agent
4
5     Functions you should fill in:
6     - computeValueFromQValues
7     - computeActionFromQValues
8     - getQValue
9     - getAction
10    - update
11
12    Instance variables you have access to
13    - self.epsilon (exploration prob)
14    - self.alpha (learning rate)
15    - self.discount (discount rate)
16
17    Functions you should use
18    - self.getLegalActions(state)
19      which returns legal actions for a state
20    """
21    def __init__(self, **args):
22        "You can initialize Q-values here..."
23        ReinforcementAgent.__init__(self, **args)
24        self.qValues = util.Counter()
25
26    def getQValue(self, state, action):
27        """
28        Returns Q(state,action)
29        Should return 0.0 if we have never seen a state
30        or the Q node value otherwise
31        """
32        return self.qValues[(state, action)] if (state, action) in self.qValues else 0.0
33
34
35    def computeValueFromQValues(self, state):
36        """
37        Returns max_action Q(state,action)
38        where the max is over legal actions. Note that if
39        there are no legal actions, which is the case at the
40        terminal state, you should return a value of 0.0.
41        """
42        legalActions = self.getLegalActions(state)
43        if not legalActions:
44            return 0.0
45        return max(self.getQValue(state, action) for action in legalActions)
46
47    def computeActionFromQValues(self, state):
48        """
49        Compute the best action to take in a state. Note that if there
50        are no legal actions, which is the case at the terminal state,
51        you should return None.
52        """
53        legalActions = self.getLegalActions(state)
54        if not legalActions:
55            return None
56
```



```

57     best_value = float('-inf')
58     best_actions = []
59
60     for action in legalActions:
61         q_value = self.getQValue(state, action)
62         if q_value > best_value:
63             best_value = q_value
64             best_actions = [action]
65         elif q_value == best_value:
66             best_actions.append(action)
67
68     return random.choice(best_actions)
69
70 def getAction(self, state):
71     """
72     Compute the action to take in the current state. With
73     probability self.epsilon, we should take a random action and
74     take the best policy action otherwise. Note that if there are
75     no legal actions, which is the case at the terminal state, you
76     should choose None as the action.
77
78     HINT: You might want to use util.flipCoin(prob)
79     HINT: To pick randomly from a list, use random.choice(list)
80     """
81     # Pick Action
82     legalActions = self.getLegalActions(state)
83     action = None
84
85     if util.flipCoin(self.epsilon):
86         action = random.choice(legalActions)
87     else:
88         action = self.computeActionFromQValues(state)
89
90     return action
91
92 def update(self, state, action, nextState, reward):
93     """
94     The parent class calls this to observe a
95     state = action => nextState and reward transition.
96     You should do your Q-Value update here
97
98     NOTE: You should never call this function,
99     it will be called on your behalf
100     """
101     sample = reward + self.discount * self.computeValueFromQValues(nextState)
102     self.qValues[(state, action)] = (1 - self.alpha) * self.getQValue(state, action) + self.
103     alpha * sample
104
105 def getPolicy(self, state):
106     return self.computeActionFromQValues(state)
107
108 def getValue(self, state):
109     return self.computeValueFromQValues(state)

```

Código 2.3: Q-Learning

Comentarios

3. Resultados