We start by instaling the MDP library that we will use throughout the workshop.

```
1   !pip install mdp
```

> Requirement already satisfied: mdp in /usr/local/lib/python3.7/dist-packages (
  Requirement already satisfied: future in /usr/local/lib/python3.7/dist-package
  Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages

```
1   import os, inspect
2   stdlib = os.path.dirname(inspect.getfile(os))
3   print(stdlib)
```

```
/usr/lib/python3.7
```

And we now mount our Drive to make use of the additional files. Here I have put these on the Colab Notebooks directoy on my Drive. If you use a different directory, the path should be changed accordingly.

```
1   from google.colab import drive
2   import sys
3   drive.mount('/content/drive')
4   sys.path.insert(0,'/content/drive/MyDrive/Colab Notebooks')
5   import _mdp as mdp
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call
```

We will now setup the practical scrpt and the virtual display for the visualisation tasks throughout. The cell may have to be run a couple of times so as to allow xvfb to start.

```
1   import sys, os
2   if 'google.colab' in sys.modules and not os.path.exists('.setup_complete'):
3     # Insert the directoryimport sys
4     # Run the scripts
5     !setup_colab_practical3.sh -O- | bash
6     !touch .setup_complete
7
8   # This code creates a virtual display to draw game images on.
9   # It will have no effect if your machine has a monitor.
10  if type(os.environ.get("DISPLAY")) is not str or len(os.environ.get("DISPLAY")
11      !bash ../xvfb start
12      os.environ['DISPLAY'] = ':1'
```

```
1   transition_probs = {
2       's0': {
3           'a0': {'s0': 0.5, 's2': 0.5},
4           'a1': {'s2': 1}
5       },
6       's1': {
```

```
 7            'a0': {'s0': 0.7, 's1': 0.1, 's2': 0.2},
 8            'a1': {'s1': 0.95, 's2': 0.05}
 9        },
10        's2': {
11            'a0': {'s0': 0.4, 's2': 0.6},
12            'a1': {'s0': 0.3, 's1': 0.3, 's2': 0.4}
13        }
14    }
15    rewards = {
16        's1': {'a0': {'s0': +5}},
17        's2': {'a1': {'s0': -1}}
18    }
19    from _mdp import MDP
20    mdp = MDP(transition_probs, rewards, initial_state='s0')
```

We can now use MDP just as any other gym environment:

```
1    print('initial state =', mdp.reset())
2    next_state, reward, done, info = mdp.step('a1')
3    print('next_state = %s, reward = %s, done = %s' % (next_state, reward, done))

     initial state = s0
     next_state = s2, reward = 0.0, done = False
```

but it also has other methods that you'll need for Value Iteration

```
1    print("mdp.get_all_states =", mdp.get_all_states())
2    print("mdp.get_possible_actions('s1') = ", mdp.get_possible_actions('s1'))
3    print("mdp.get_next_states('s1', 'a0') = ", mdp.get_next_states('s1', 'a0'))
4    print("mdp.get_reward('s1', 'a0', 's0') = ", mdp.get_reward('s1', 'a0', 's0'))
5    print("mdp.get_transition_prob('s1', 'a0', 's0') = ", mdp.get_transition_prob(

     mdp.get_all_states = ('s0', 's1', 's2')
     mdp.get_possible_actions('s1') =  ('a0', 'a1')
     mdp.get_next_states('s1', 'a0') =  {'s0': 0.7, 's1': 0.1, 's2': 0.2}
     mdp.get_reward('s1', 'a0', 's0') =  5
     mdp.get_transition_prob('s1', 'a0', 's0') =  0.7
```

## ▾ Visualising MDPs

You can visualize any MDP with the drawing fuction donated by [neer201](#). To do so, we have to install graphviz.

```
1    from _mdp import has_graphviz
2    from IPython.display import display
3    print("Graphviz available:", has_graphviz)

     Graphviz available: True
```
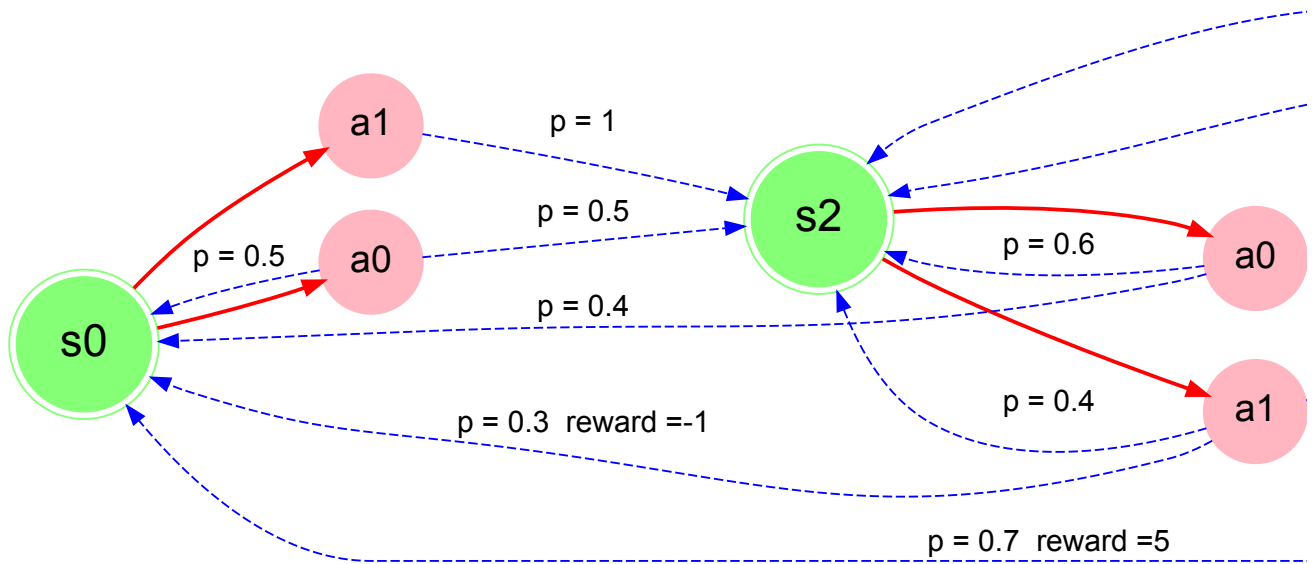
```
1    if has graphviz:
```

```
2        from _mdp import plot_graph, plot_graph_with_state_values, plot_graph_opti
3        display(plot_graph(mdp))
```



## Policy iteration algorithm

Here's the pseudo-code for PI:

1. Pick an random policy $\pi$

2. Iterate until $\pi$ is unchanged (converges).

**Policy evalution:**

$$V_{i+1}(s) = \sum_{s'} P(s'|s, \pi(s)) \cdot [r(s, \pi(s), s') + \gamma V_i(s')]$$

until value converges

**Policy improvement:**

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

First, let's write a function to compute the state-action value function $Q(\pi)$, defined as follows

$$Q_i(s, a) = \sum_{s'} P(s'|s, a) \cdot [r(s, a, s') + \gamma V_i(s')]$$

Using $Q(s, a)$ we calculate the value at state $s$ with policy $\pi(s)$.

$$V_{i+1}(s) = \sum_{s'} P(s'|s, \pi(s)) \cdot [r(s, \pi(s), s') + \gamma V_i(s')]$$

and loop until the value if coverage, the different from current loop with previous is not much.

Then, We try update the policy base on the new value of $V$

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

loop until $\pi$ is converges.

```
1   # Compute the Q-value using the formula above
2   def get_action_value(mdp, state_values, state, action, gamma):
3       # Initialise Q
4       Q = 0
5       for s in mdp.get_all_states():
6         # Compute Q using the equation above
7         Q = Q + mdp.get_transition_prob(state, action, s)*(mdp.get_reward(state,
8                                                 gamma*state_val
9       return Q
```

```
1   # initial simple policy
2   from collections import defaultdict
3   def init_policy():
4       policy = defaultdict(lambda: {})
5       for state in mdp.get_all_states():
6           actions = mdp.get_possible_actions(state)
7           for action in actions:
8               policy[state][action] = 1. / len(actions)
9       return policy
```

```
1   def policy_eval(mdp, policy, gamma, theta=0.00001):
2       num_iter = 100  # maximum iterations, excluding initialization
3       # initialize V(s)
4       V = {s: 0 for s in mdp.get_all_states()}
5
6       for i in range(num_iter):
7           delta = 0
8           # Compute new state values using the functions defined above.
9           new_V = {}
10          for s in mdp.get_all_states():
11            nsv = get_action_table(mdp, V, s, gamma)
12            # new_V[s] = max(nsv.items(), key=operator.itemgetter(1))[1]
13            action_policy = best_action(policy[s])
14            new_V[s] = nsv[action_policy]
15
16            delta = max(delta, abs(new_V[s] - V[s]))
17          assert isinstance(new_V, dict)
18          V = new_V
19          if delta < theta:
20              break
21
22      return V
```

```
1   def get_action_table(mdp, state_values, state, gamma):
2       A = {a: 0 for a in mdp.get_possible_actions(state)}
3       i = 0
4       # Compute all possible options
5       for a in mdp.get_possible_actions(state):
6           v = get_action_value(mdp, state_values, state, a, gamma)
7           A[a] = v
8       return A
```

```
1   def best_action(policy_state):
2     return max(policy_state.items(), key=operator.itemgetter(1))[0] # argmax
```

```
1   import operator
2   def policy_improvement(mdp, policy_eval_fn=policy_eval, gamma=0.9):
3       num_iter = 100  # maximum iterations, excluding initialization
4
5       policy = init_policy()
6       for i in range(num_iter):
7         V = policy_eval_fn(mdp, policy, gamma)
8         print(V)
9         # Will be set to false if we make any changes to the policy
10        policy_stable = True
11        # For each state...
12        for state in mdp.get_all_states():
13            # The best action we would take under the current policy
14            chosen_a = best_action(policy[state])
15
16            action_values = get_action_table(mdp, V, state, gamma)
17            best_a = best_action(action_values)
18            # Greedily update the policy
19            if chosen_a != best_a:
20                policy_stable = False
21            for action in policy[state]:
22              if action == best_a:
23                policy[state][action] = 1
24              else:
25                policy[state][action] = 0
26        # If the policy is stable we've found an optimal policy. Return it
27        if policy_stable:
28            return policy, V
29      return policy, V
```

```
1   import pprint
2   policy, v = policy_improvement(mdp)
3   print("Policy Probability Distribution:")
4   pprint.pprint(policy)
5   print("")
6   print(v)
```

```
{'s0': 0.0, 's1': 3.8461536621935, 's2': 0.0}
{'s0': 2.8398740312453152, 's1': 6.498777111635999, 's2': 3.4709784640026404}
{'s0': 3.7898252943341277, 's1': 7.302796844653717, 's2': 4.210930696013536}
Policy Probability Distribution:
defaultdict(<function init_policy.<locals>.<lambda> at 0x7f651d5efb90>,
            {'s0': {'a0': 0, 'a1': 1},
             's1': {'a0': 1, 'a1': 0},
             's2': {'a0': 0, 'a1': 1}})

{'s0': 3.7898252943341277, 's1': 7.302796844653717, 's2': 4.210930696013536}
```
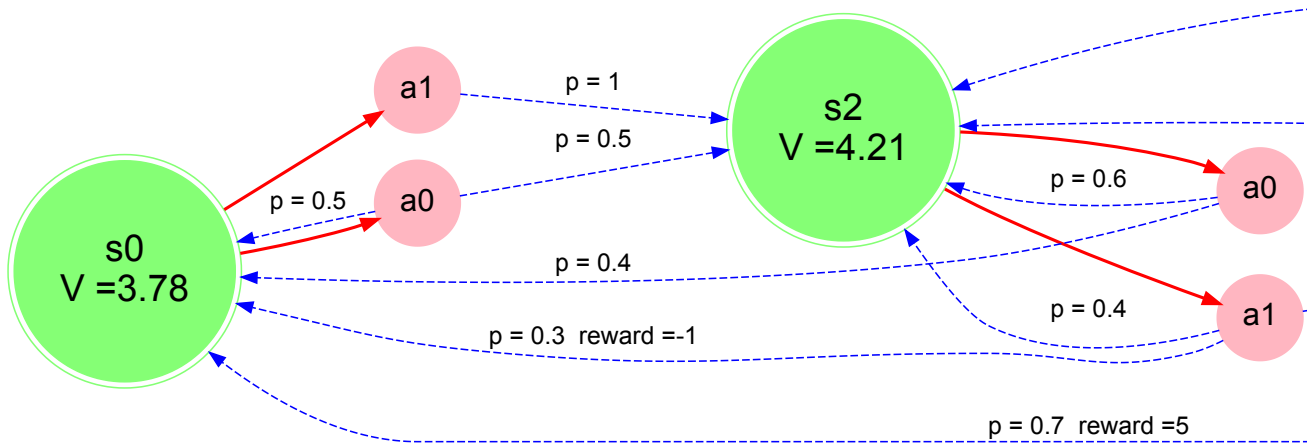
Finally, we can now plot the final graphical model

```
1  if has_graphviz:
2      display(plot_graph_with_state_values(mdp, v))
```
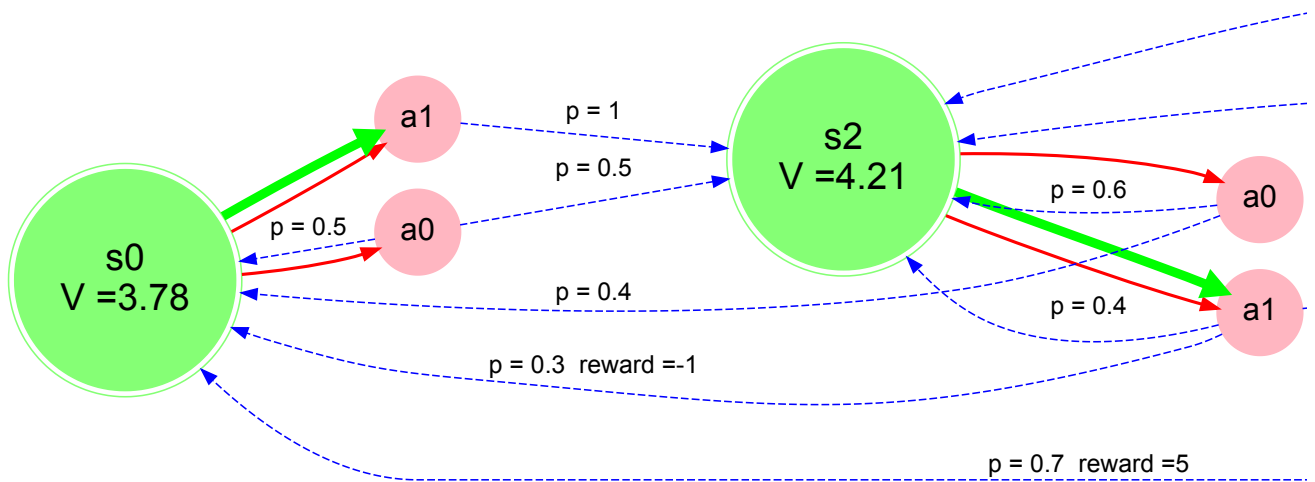


Also, we can now look at te final state values.

```
1  if has_graphviz:
2      display(plot_graph_optimal_strategy_and_state_values(mdp, v, get_action_va
```



```
1  def get_optimal_action(mdp, state):
2      # Finds optimal action using formula above.
3      if mdp.is_terminal(state):
4          return None
5
6      return best_action(policy[state])
```

```
1  gamma=0.9
2  assert get_optimal_action(mdp, 's0') == 'a1'
3  assert get_optimal_action(mdp, 's1') == 'a0'
4  assert get_optimal_action(mdp, 's2') == 'a1'
```

```
1  import numpy as np
```

```
2    # Measure agent's average reward
3
4    s = mdp.reset()
5    rewards = []
6    gamma=0.9
7
8    for _ in range(1000):
9        s, r, done, _ = mdp.step(get_optimal_action(mdp, s))
10       rewards.append(r)
11
12   print("average reward: ", np.mean(rewards))
13
14   assert(0.40 < np.mean(rewards) < 0.55)
```

```
average reward:  0.453
```

1