In [8]:
```python
import numpy as np
from enum import Enum
import copy
```
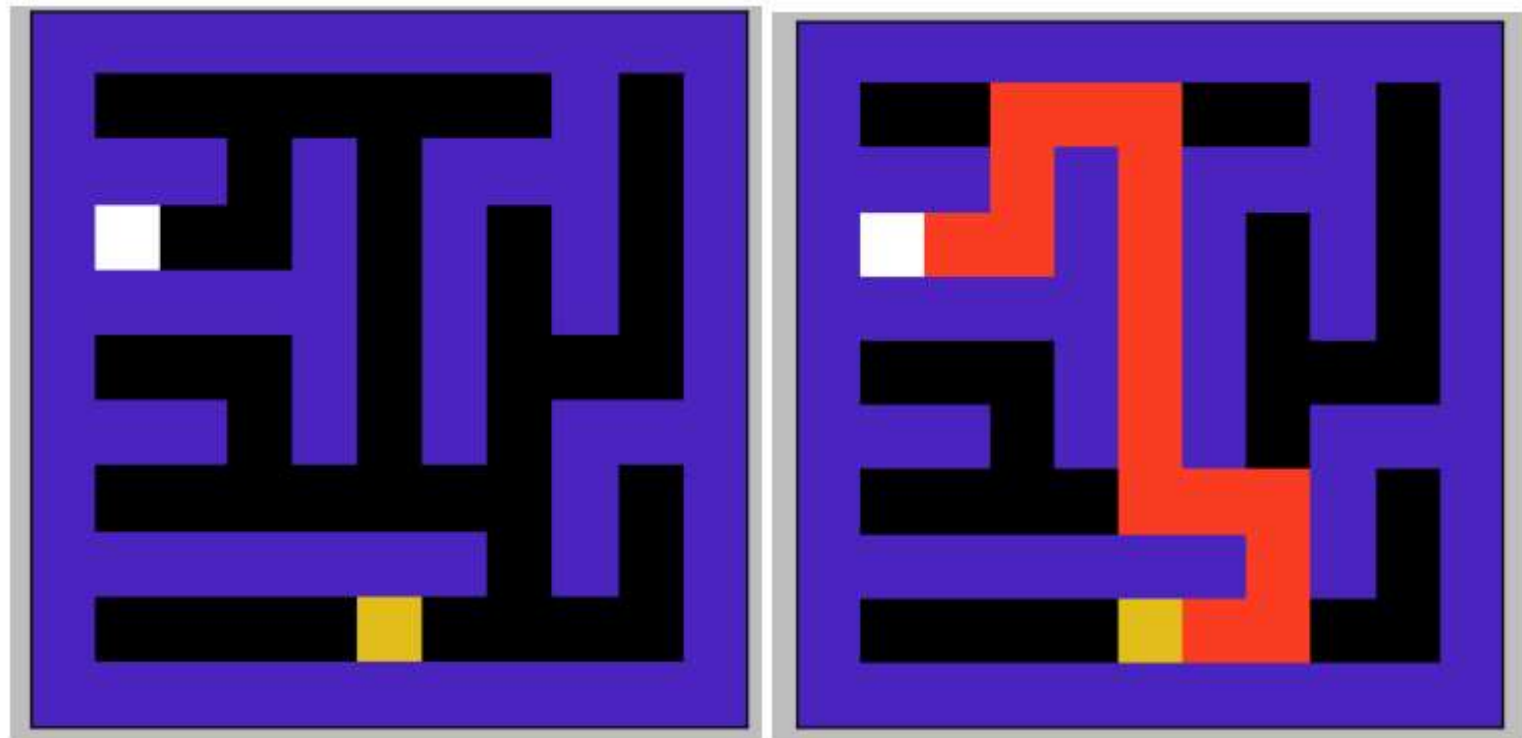
Consider a standard grid world, where only 4 (up, down, left, right) actions are allowed and the agent deterministically moves accordingly, represented as below. Here yellow is the start state and white is the goal state.

Say, we define our MDP as:

- S: 121 (11 x 11) cells
- A: 4 actions (up, down, left, right)
- P: Deterministic transition probability
- R: -1 at every step
- gamma: 0.9

Our goal is to find an optimal policy (shown in right).



In [9]:
```python
# Above grid is defined as below:
#   - 0 denotes an navigable tile
#   - 1 denotes an obstruction/wall
#   - 2 denotes the start state
#   - 3 denotes an goal state

# Note: Here the upper left corner is defined as (0, 0)
#       and lower right corner as (m-1, n-1)

# Optimal Path: RIGHT RIGHT UP UP LEFT LEFT UP UP UP UP UP UP LEFT LEFT DOWN DOWN LEFT LEFT


GRID_WORLD = np.array([
    [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],
    [1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1],
    [1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1],
    [1, 3, 0, 0, 1, 0, 1, 0, 1, 0, 1],
    [1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1],
    [1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1],
    [1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1],
    [1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1],
    [1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1],
    [1, 0, 0, 0, 0, 2, 0, 0, 0, 0, 1],
    [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
])
```

## Actions

```
In [10]:   1  class Actions(Enum):
           2      UP    = (0, (-1, 0))  # index = 0, (xaxis_move = -1 and yaxis_move = 0)
           3      DOWN  = (1, (1, 0))   # index = 1, (xaxis_move = 1 and yaxis_move = 0)
           4      LEFT  = (2, (0, -1))  # index = 2, (xaxis_move = 0 and yaxis_move = -1)
           5      RIGHT = (3, (0, 1))   # index = 3, (xaxis_move = 0 and yaxis_move = -1)
           6
           7      def get_action_dir(self):
           8          _, direction = self.value
           9          return direction
          10
          11      @property
          12      def index(self):
          13          indx, _ = self.value
          14          return indx
          15
          16      @classmethod
          17      def from_index(cls, index):
          18          action_index_map = {a.index: a for a in cls}
          19          return action_index_map[index]
```

```
In [11]:   1  # How to use Action enum
           2  for a in Actions:
           3      print(f"name: {a.name}, action_id: {a.index}, direction_to_move: {a.get_action_dir()}")
           4
           5  print("\n-----------------------------------\n")
           6
           7  # find action enum from index 0
           8  a = Actions.from_index(3)
           9  print(f"3 index action is: {a.name}")
```

```
name: UP, action_id: 0, direction_to_move: (-1, 0)
name: DOWN, action_id: 1, direction_to_move: (1, 0)
name: LEFT, action_id: 2, direction_to_move: (0, -1)
name: RIGHT, action_id: 3, direction_to_move: (0, 1)

-----------------------------------

3 index action is: RIGHT
```

## Policy

```
In [12]:   1  class BasePolicy:
           2      def update(self, *args):
           3          pass
           4
           5      def select_action(self, state_id: int) -> int:
           6          raise NotImplemented
           7
           8
           9  class DeterministicPolicy(BasePolicy):
          10      def __init__(self, actions: np.ndarray):
          11          # actions: its a 1d array (|S| size) which contains action for each state
          12          self.actions = actions
          13
          14      def update(self, state_id, action_id):
          15          assert state_id < len(self.actions), f"Invalid state_id {state_id}"
          16          assert action_id < len(Actions), f"Invalid action_id {action_id}"
          17          self.actions[state_id] = action_id
          18
          19      def select_action(self, state_id: int) -> int:
          20          assert state_id < len(self.actions), f"Invalid state_id {state_id}"
          21          return self.actions[state_id]
```

## Environment

In [13]:
```python
class Environment:
  def __init__(self, grid):
    self.grid = grid
    m, n = grid.shape
    self.num_states = m*n

  def xy_to_posid(self, x: int, y: int):
    _, n = self.grid.shape
    return x*n + y

  def posid_to_xy(self, posid: int):
    _, n = self.grid.shape
    return (posid // n, posid % n)

  def isvalid_move(self, x: int, y: int):
    m, n = self.grid.shape
    return (x >= 0) and (y >= 0) and (x < m) and (y < n) and (self.grid[x, y] != 1)

  def find_start_xy(self) -> int:
    m, n = self.grid.shape
    for x in range(m):
      for y in range(n):
        if self.grid[x, y] == 2:
          return (x, y)
    raise Exception("Start position not found.")

  def find_path(self, policy: BasePolicy) -> str:
    max_steps = 50
    steps = 0

    P, R = self.get_transition_prob_and_expected_reward()
    num_actions, num_states = R.shape
    all_possible_state_posids = np.arange(num_states)

    path = ""
    curr_x, curr_y = self.find_start_xy()
    while (self.grid[curr_x, curr_y] != 3) and (steps < max_steps):
      curr_posid = self.xy_to_posid(curr_x, curr_y)
      action_id = policy.select_action(curr_posid)
      next_posid = np.random.choice(
          all_possible_state_posids, p=P[action_id, curr_posid])
      action = Actions.from_index(action_id)
      path += f" {action.name}"
      curr_x, curr_y = self.posid_to_xy(next_posid)
      steps += 1
    return path

  def get_transition_prob_and_expected_reward(self):  # P(s_next | s, a), R(s, a)
    m, n = self.grid.shape
    num_states = m*n
    num_actions = len(Actions)
    P = np.zeros((num_actions, num_states, num_states))
    R = np.zeros((num_actions, num_states))
    for a in Actions:
      for x in range(m):
        for y in range(n):
          xmove_dir, ymove_dir = a.get_action_dir()
          xnew, ynew = x + xmove_dir, y + ymove_dir  # find the new co-ordinate after the action a

          posid = self.xy_to_posid(x, y)
          new_posid = self.xy_to_posid(xnew, ynew)


          if self.grid[x, y] == 3:
            # the current state is a goal state
            P[a.index, posid, posid] = 1
            R[a.index, posid] = 0
          elif (self.grid[x, y] == 1) or (not self.isvalid_move(xnew, ynew)):
            # the current state is a block state or the next state is invalid
            P[a.index, posid, posid] = 1
            R[a.index, posid] = -1
          else:
            # action a is valid and goes to a new position
            P[a.index, posid, new_posid] = 1
            R[a.index, posid] = -1
    return P, R
```

## Policy Iteration

In [14]:

```python
def policy_evaluation(P: np.ndarray, R: np.ndarray, gamma: float,
                      policy: BasePolicy, theta: float,
                      init_V: np.ndarray=None):
    num_actions, num_states = R.shape

    # Please try different starting point for V you will find it will always
    # converge to the same V_pi value.
    if init_V is None:
        init_V = -2.5*np.ones(num_states) # Different Starting point for V
    V = copy.deepcopy(init_V)

    delta = 100.0
    while delta > theta:
        delta = 0.0
        for state_id in range(num_states):
            action_id = policy.select_action(state_id)
            v_old = V[state_id]
            # Following equation is a different way of writing the same equation given in the slide.
            # Note here R is an expected reward term.
            V[state_id] = R[action_id, state_id] + gamma * np.dot(P[action_id, state_id], V)
            delta = max(delta, abs(V[state_id] - v_old))
    return V


def policy_improvement(P: np.ndarray, R: np.ndarray, gamma: float,
                       policy: BasePolicy, V: np.ndarray):
    num_actions, num_states = R.shape
    policy_stable = True
    for state_id in range(num_states):
        old_action_id = policy.select_action(state_id)

        # your code here
        r = R[:, state_id]
        p = P[:, state_id]
        new_action_id = np.argmax(r + gamma * np.dot(p, V)) # update new_action_id based on the value function.

        policy.update(state_id, new_action_id)
        if old_action_id != new_action_id:
            policy_stable = False
    return policy_stable


def policy_iteration(P: np.ndarray, R: np.ndarray, gamma: float,
                     theta: float=1e-3, init_policy: BasePolicy = None):
    num_actions, num_states = R.shape

    # Please try exploring different policies you will find it will always
    # converge to the same optimal policy for valid states.
    if init_policy is None:
        # Say initial policy = all down actions.
        init_policy = DeterministicPolicy(actions=np.ones(num_states, dtype=int)) # Changed All initial actions to go

    # creating a copy of a initial policy
    policy = copy.deepcopy(init_policy)
    policy_stable = False
    while not policy_stable:
        V = policy_evaluation(P, R, gamma, policy, theta)
        policy_stable = policy_improvement(P, R, gamma, policy, V)
    return policy, V
```

## Value Iteration

```python
In [15]:   1  # Directly find the optimal value function
           2  def get_optimal_value(P: np.ndarray, R: np.ndarray, gamma: float,
           3                        theta: float, init_V: np.ndarray=None):
           4    num_actions, num_states = R.shape
           5
           6    # Please try different starting point for V you will find it will always
           7    # converge to the same V_star value.
           8    if init_V is None:
           9      init_V = -6.75* np.ones(num_states) # Different Starting point for V
          10    V = copy.deepcopy(init_V)
          11
          12    delta = 100.0
          13    while delta > theta:
          14      delta = 0.0
          15      for state_id in range(num_states):
          16        v_old = V[state_id]
          17        q_sa = np.zeros(num_actions)
          18        for a in Actions:
          19          q_sa[a.index] = R[a.index, state_id] + gamma * np.dot(P[a.index, state_id], V)
          20        V[state_id] = np.max(q_sa)
          21        delta = max(delta, abs(V[state_id] - v_old))
          22    return V
          23
          24
          25  def value_iteration(P: np.ndarray, R: np.ndarray, gamma: float,
          26                      theta: float=1e-3, init_V: np.ndarray=None):
          27    V_star = get_optimal_value(P, R, gamma, theta, init_V)
          28
          29    num_actions, num_states = R.shape
          30    policy = DeterministicPolicy(actions=np.zeros(num_states, dtype=int))
          31    for state_id in range(num_states):
          32      # Your code here
          33      r = R[:, state_id]
          34      p = P[:, state_id]
          35      action_id = np.argmax(r + gamma * np.dot(p, V_star)) # update the action_id based on V_star
          36
          37      policy.update(state_id, action_id)
          38
          39    return policy, V_star
```

## Experiments

```python
In [16]:   1  def is_same_optimal_value(V1, V2, diff_theta=1e-3):
           2    diff = np.abs(V1 - V2)
           3    return np.all(diff < diff_theta)
```

```python
In [17]:   1  seed = 0
           2  np.random.seed(seed)
           3
           4  gamma = 0.9
           5  theta = 1e-5
```

```python
In [18]:   1  env = Environment(GRID_WORLD)
           2  P, R = env.get_transition_prob_and_expected_reward()
```

**Exp 1: Using Policy iteration algorithm find the optimal path from start to goal position**

```python
In [19]:   1  # # Start with random choice of init_policy.
           2  # One such choice could be: init_policy = np.ones(env.num_states, dtype=int)
           3  # Start with your own choice of init_policy
           4  init_policy = DeterministicPolicy(actions=2 * np.ones(env.num_states, dtype=int)) # Intial actions set to 2: LEFT
           5
           6  pitr_policy, pitr_V_star = policy_iteration(P, R, gamma, theta=theta, init_policy=init_policy)
           7  pitr_path = env.find_path(pitr_policy)
           8  print(pitr_path)
```

```
RIGHT RIGHT UP UP LEFT LEFT UP UP UP UP UP UP LEFT LEFT DOWN DOWN LEFT LEFT
```

**Exp 2: Using value iteration algorithm find the optimal path from start to goal position**

```python
In [20]:    1  vitr_policy, vitr_V_star = value_iteration(P, R, gamma, theta=theta)
            2  vitr_path = env.find_path(vitr_policy)
            3  print(vitr_path)
```

```
RIGHT RIGHT UP UP LEFT LEFT UP UP UP UP UP UP LEFT LEFT DOWN DOWN LEFT LEFT
```

**Exp 3: Compare the optimal value function of policy iteration and value iteration algorithm**

```python
In [22]:    1  is_same_optimal_value(pitr_V_star, vitr_V_star)
```

Out[22]:  True

**Exp 4: Using initial guess for V as random values, find the optimal value function using policy evaluation and compare it with the optimal value function**

```python
In [23]:    1  # Start with random choice of init_V.
            2  # One such choice could be: init_V = np.random.randn(env.num_states)
            3  # Another choice could be: init_V = 10*np.ones(env.num_states)
            4  # Start with your own choice of init_V
            5  init_V = -169*np.random.rand(env.num_states) # Different Starting point for V
            6
            7  V_star = policy_evaluation(P, R, gamma, pitr_policy, theta, init_V)
            8  is_same_optimal_value(pitr_V_star, V_star)
```

Out[23]:  True

**Exp 5: Using initial guess for V as random values, find the optimal value function using get_optimal_value and compare it with the optimal value function**

```python
In [24]:    1  # Start with random choice.
            2  # One such choice could be: init_V = np.random.randn(env.num_states)
            3  # Another choice could be: init_V = 10*np.ones(env.num_states)
            4  # Start with your own choice of init_V
            5  init_V = -1e7*np.random.random(env.num_states) # Different Starting point for V
            6
            7  V_star = get_optimal_value(P, R, gamma, theta, init_V)
            8  is_same_optimal_value(vitr_V_star, V_star)
```

Out[24]:  True

**Exp Optional: Try changing the grid by adding multiple paths to the goal state and check if our policy_iteration or value_iteration algorithm is able to find optimal path. Redo the above experiments.**

- 1 way to add another path would be GRID_WORLD[4, 1] = 0

```python
In [ ]:    1
```

```python
In [25]:    1  from google.colab import drive
            2  drive.mount('/content/drive')
```

```
Mounted at /content/drive
```

```python
In [28]:    1  !jupyter nbconvert --to pdf /content/drive/MyDrive/DA6400_Value_and_Policy_Iteration_DA24S018.ipynb
```

```
[NbConvertApp] Converting notebook /content/drive/MyDrive/DA6400_Value_and_Policy_Iteration_DA24S018.ipynb to pdf
[NbConvertApp] Writing 63602 bytes to notebook.tex
[NbConvertApp] Building PDF
[NbConvertApp] Running xelatex 3 times: ['xelatex', 'notebook.tex', '-quiet']
[NbConvertApp] Running bibtex 1 time: ['bibtex', 'notebook']
[NbConvertApp] WARNING | bibtex had problems, most likely because there were no citations
[NbConvertApp] PDF successfully created
[NbConvertApp] Writing 68870 bytes to /content/drive/MyDrive/DA6400_Value_and_Policy_Iteration_DA24S018.pdf
```

```python
In [ ]:    1
```