

這段程式碼展示了一個簡單的網格世界環境（**GridWorld Environment**），用於強化學習中的決策問題。以下是主要步驟和意義：

1. ****環境設定****：

- 程式碼定義了一個 `GridWorldEnv` 類，繼承自 `gym.Env`，用於創建和管理一個 4x6 的網格世界環境。每個單元格大小為 100 像素，並設有一些預設的參數，如延遲時間（`delay`）。

2. ****狀態與動作****：

- 環境中有 24 個狀態（每個網格單元格代表一個狀態），並且有 4 個動作（上、下、左、右）。這些動作改變代理人的位置。

3. ****終端狀態****：

- 定義了金幣位置（目標狀態）和陷阱位置（負面狀態）。代理人達到金幣位置時獲得正獎勵，達到陷阱位置時獲得負獎勵。

4. ****轉移概率****：

- 設置了每個狀態和動作對應的轉移概率。根據當前狀態和動作，計算轉移到下一個狀態的獎勵和是否完成。

5. ****顯示界面****：

- 使用 `matplotlib` 繪製環境的圖形界面，包括網格、陷阱和金幣位置。代理人的位置會根據 `pkl` 文件中的記錄來顯示，或者預設在 (0,0) 位置。

6. ****重置與步驟****：

- `reset()` 方法將環境重置到初始狀態。`step(action)` 方法根據代理人的

動作更新狀態，返回下一個狀態、獎勵和是否達到終端狀態。

7. ****渲染與關閉****：

- ``render(mode='human', done=False)`` 方法用於渲染當前狀態，使代理人的位置可視化。``close()`` 方法關閉圖形界面。

8. ****主程序****：

- 在主程序中，初始化環境並隨機選擇動作進行測試，顯示每個步驟的結果，直到達到終端狀態為止。

這段程式碼提供了一個基於網格的強化學習環境的基本實現，可以用於測試和開發強化學習算法。