## Q1.

**(1,2).** The solutions of part1 and part2 are given by excel r=-0.05 and r=-0.1

**(3).** The state for which the optimal policies in iteration 10 is different for the two values of R(s) is U9(12). The policy is left when r=-0.05. But the policy is down when r=-0.1.

When r=-0.05, the policy is left. Because the cost of taking a step is fairly small compared with the penalty for ending up in (4,2) by accident, the optimal policy for the state (3,1) is conservative. The policy recommends taking the long way round, rather than taking the shortcut and thereby risking entering (4,2).

When r=-0.1, the policy is left. This is because the discount number is different. We can intuitively interpret this as the when R(s)=-0.1 the penalty is serious and life is quite unpleasant. So the agent takes the shortest route to the +1 state and is willing to risk falling into the −1 state by accident.

**(4).** The optimal policies for $R(s) = -0.04$ and $R(s) = -0.05$ are not identical. The policies for $R(s) = -0.04$ given by following.

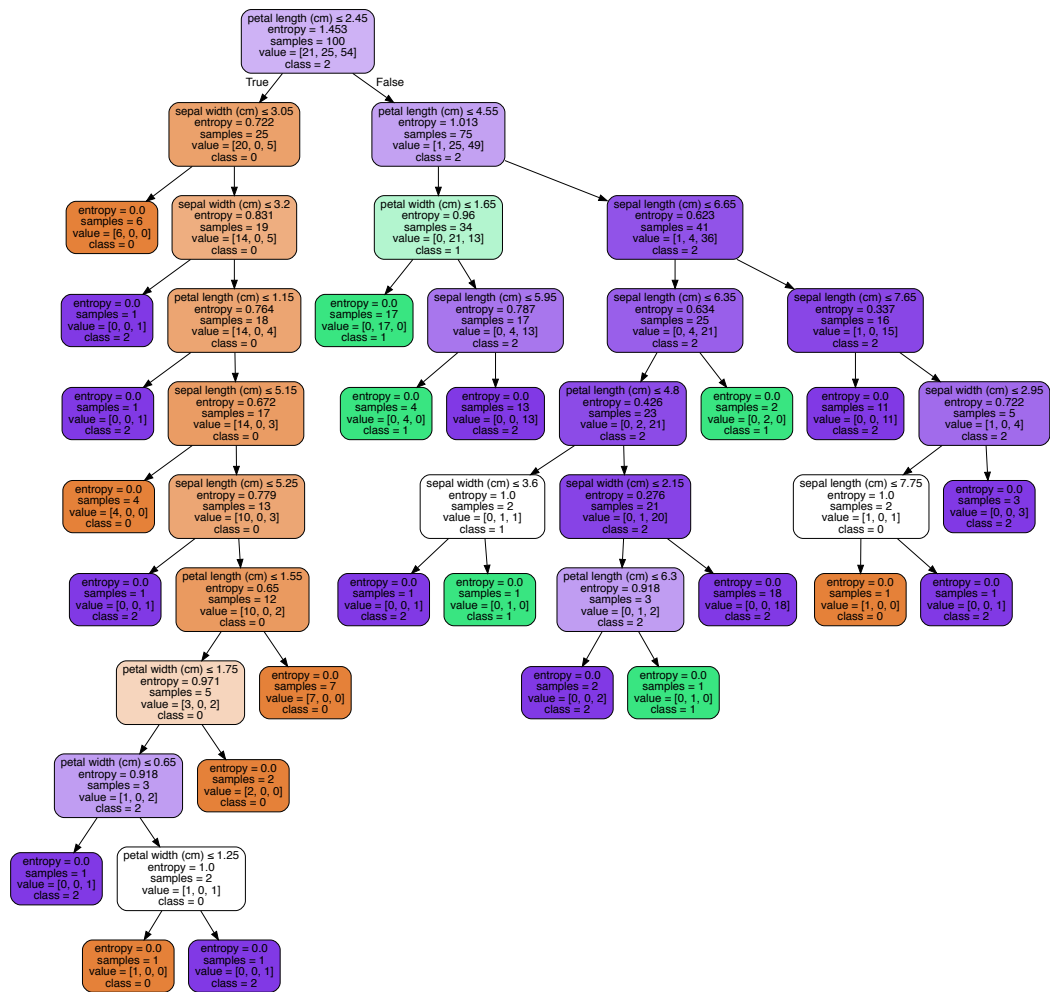|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | ↓ | ← | ← | ← |
| 2 | ↓ | X | ↓ |  -1 |
| 3 | → | → | → | 1 |

The policies for $R(s) = -0.05$ is given by follows.

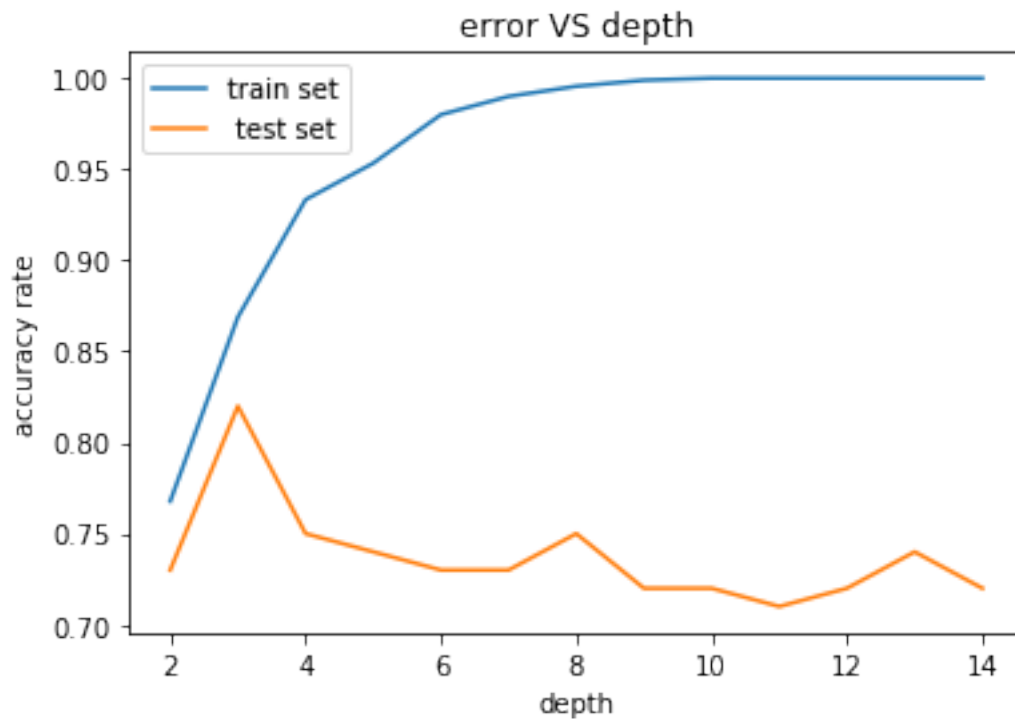|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | ↓ | ← | ↓ | ← |
| 2 | ↓ | X | ↓ | -1 |
| 3 | → | → | → | 1 |

The row1col3 is different. I think this is because when $R(s) = -0.05$, the life is more painful than that of $R(s) = -0.04$, the agent takes the shortest route to the +1 state.
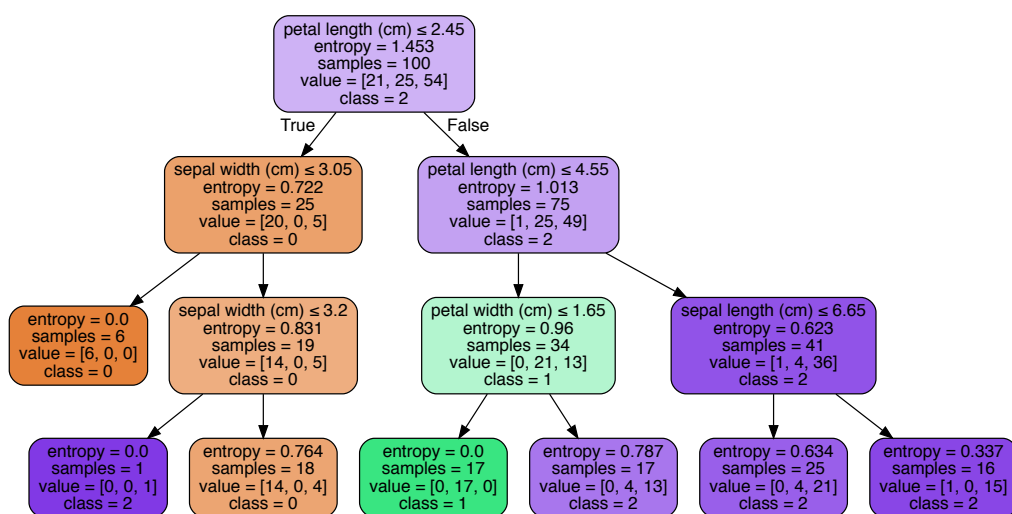
**Q2**

**1.**

**2.(2)**



**(3)**

As shown on above graph when depth=3 the accuracy of validation is highest.

**(4)**

**(5)**

The prediction accuracy of train set with depth=3 is 0.933

The prediction accuracy of test set with depth=3 is 0.8  0.869

The prediction accuracy of original set with depth=3 is 0.87


(The code for question 2.1, 2.2 and 2.3 are in the folder.)