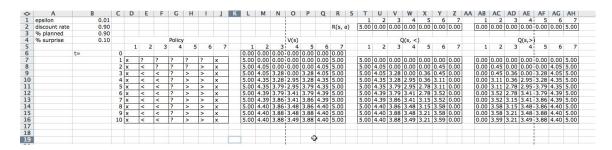# Reinforcement Learning Self Check
# 605.645 – Artificial Intelligence

The purpose of this self-check is to make sure you understand key concepts for the algorithms presented during the module and to prepare you for the programming assignment. As you work through problems, you should always be thinking "how would I do this in code? What basic data structures would I need? What operations on those basic data structures?"

As noted in the video lectures, it is possible to solve value iteration for a 1-d Markov Decision Process in a spreadsheet by explicitly modeling the π, V, and Q arrays. You have been provided with a starter spreadsheet (Module 11 Self Check.xlsx)

For the self-check, you have been provided with a starter spreadsheet (module4-self-check.xlsx). Fill in the formulas for Value Iteration (stochastic version) so that it looks like the following:



1.  the cell B1 is named "discount_rate", you can refer to it by that name in cell formulas, for example, "=discount_rate*10"
2.  the cell B2 is named "planned", you can refer to it by that name in formulas.
3.  The cell B3 is named "surprise", you can refer to it by that name in formulas.
4.  if(condition,then,else) can be used to test values in cells. It can also be nested (the first one will need an =).
5.  =max(*cellref1*, *cellref2*) will return the maximum value.
6.  You do not need to fill in formulas for the goal states 1 or 7 except in the case of V(s).
7.  Use "<" for "left", "?" for "pick random" and ">" for "right" in your policy formula.
8.  There isn't a good way to control the number of iterations automatically, so epsilon is there to let you know when you should stop copying.

You should get the same results as above and converge in 10 iterations.

**Additional**

1. Copy everything into a new Tab, "Right". Change the rewards so that the agent always goes right.
2. Copy everything into a new Tab, "Left". Change the rewards so that the agent always goes left except if they're in the square right next to the goal on the right.
3. Try make a new Tab, "Fourteen", that increases the size of the world from 7 to 14. Leave the rewards the same values. What happens to the policy? Use your experiments with #1 and #2 to change the rewards so that there is a deterministic policy for the whole world (no "?").