



[PYTHON \(https://www.analyticsvidhya.com/blog/category/python-2/\)](https://www.analyticsvidhya.com/blog/category/python-2/)

[REINFORCEMENT LEARNING \(https://www.analyticsvidhya.com/blog/category/machine-learning/reinforcement-learning/\)](https://www.analyticsvidhya.com/blog/category/machine-learning/reinforcement-learning/)

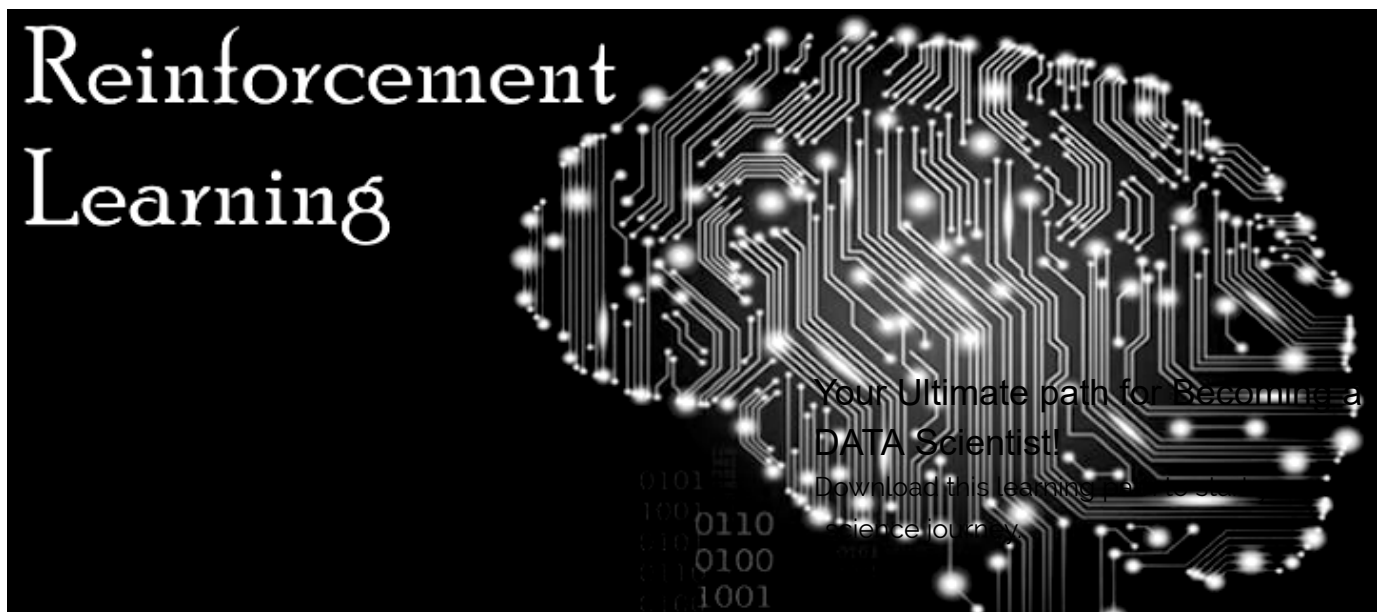
## Nuts & Bolts of Reinforcement Learning: Model Based Planning using Dynamic Programming

ANKIT CHOUDHARY (<https://www.analyticsvidhya.com/blog/author/ankit2106/>), SEPTEMBER 18, 2018

[LOGIN TO BOOKMARK THIS ARTICLE \(https://id.analyticsvidhya.com/accounts/login/?next=https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/\)](https://id.analyticsvidhya.com/accounts/login/?next=https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/)

### Introduction

Deep Reinforcement learning is responsible for the two biggest AI wins over human professionals – Alpha Go and OpenAI Five. Championed by Google and Elon Musk, interest in this field has gradually increased in recent years to the point where it's a thriving area of research nowadays.



In this article, however, we will not talk about a typical RL setup but explore Dynamic Programming (DP). DP is a collection of algorithms that can solve a problem *where we have the perfect model of the environment (i.e. probability distributions of any change happening in the problem setup are known) and where an agent can only take discrete actions.*

DP essentially solves a planning problem rather than a more general RL problem. The main difference, as mentioned, is that for an RL problem the environment can be very complex and its specifics are not known at all initially.

Name

Email

Contact Number



Download Resource

But before we dive into all that, let's understand why you should learn dynamic programming in the first place using an intuitive example.

## Why learn dynamic programming?

Apart from being a good starting point for grasping reinforcement learning, dynamic programming can help find optimal solutions to planning problems faced in the industry, with an important assumption that the specifics of the environment are known. DP presents a good starting point to understand RL algorithms that can solve more complex problems.

## Sunny's Motorbike Rental company

Sunny manages a motorbike rental company in Ladakh. Being near the highest motorable road in the world, there is a lot of demand for motorbikes on rent from tourists. Within the town he has 2 locations where tourists can come and get a bike on rent. If he is out of bikes at one location, then he loses business.



## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.

- Bikes are rented out for Rs 1200 per day and are available for renting the day after they are returned.
- Sunny can move the bikes from 1 location to another and incurs a cost of Rs 100.
- With experience Sunny has figured out the approximate probability distributions of demand and return rates.
- Number of bikes returned and requested at each location are given by functions  $g(n)$  and  $h(n)$  respectively. In exact terms the probability that the number of bikes rented at both locations is  $n$  is given by  $g(n)$  and probability that the number of bikes returned at both locations is  $n$  is given by  $h(n)$

Name

The problem that Sunny is trying to solve is to find out how many bikes he should move each day from 1 location to another so that he can maximise his earnings.

Email Id

Here, we exactly know the environment ( $g(n)$  &  $h(n)$ ) and this is the kind of problem in which dynamic programming can come in handy. Similarly, if you can properly model the environment of your problem where you can take discrete actions, then DP can help you find the optimal solution. In this article, we will use DP to train an agent using Python to traverse a simple environment, while touching upon key concepts in RL such as policy, reward, value function and more.

Contact Number

Download Resource

## Table of Contents

1. Understanding Agent-Environment interface using tic-tac-toe
2. Introduction to Markov Decision Process
  1. Value Function:
    1. How good it is to be in a given state?
    2. How good an action is at a particular state?
  2. Solving an MDP: The Math
    1. Bellman Expectation Equation
    2. Bellman Optimality Equation
3. Dynamic Programming
  1. Policy Iteration
    1. Policy Evaluation
    2. Policy Improvement
  2. Value Iteration
4. DP in action: Finding optimal policy for Frozen Lake environment using Python
  1. Frozen Lake Environment
  2. Policy Iteration in python
  3. Value Iteration in python

## Understanding Agent Environment Interface using tic-tac-toe

Most of you must have played the tic-tac-toe game in your childhood. If not, you can grasp the rules of this simple game from its [wiki page \(https://en.wikipedia.org/wiki/Tic-tac-toe\)](https://en.wikipedia.org/wiki/Tic-tac-toe). Suppose tic-tac-toe is your favourite game, but you have nobody to play it with. So you decide to design a bot that can play this game with you. Some key questions are:

Can you define a rule-based framework to design an efficient bot?

You sure can, but you will have to hardcode a lot of rules for each of the possible situations that might arise in a game. However, an even more interesting question to answer is:

Can you train the bot to learn by playing against you several times? And that too without being explicitly programmed to play tic-tac-toe efficiently?

A few considerations for this are:

- First, the bot needs to understand the situation it is in. A tic-tac-toe has 9 spots to fill with an X or O. Each different possible combination in the game will be a different situation for the bot, based on which it will make the next move. Each of these scenarios as shown in the below image is a different **state**.

**Your Ultimate path for Becoming a DATA Scientist!**

Download this learning path to start your data science journey.



Download Resource

	O	
O	X	

	O	
O	X	X

		X
	O	
O	X	

- Once the state is known, the bot must take an **action** in a way it considers to be optimum to win the game (**policy**)
- This move will result in a new scenario with new combinations of O's and X's which is a **new state** and a numerical **reward** will be given based on the quality of move with the goal of winning the game (**cumulative reward**)

For more clarity on the aforementioned reward, let us consider a match between bots O and X:



Consider the following situation encountered in tic-tac-toe:

	O	
O	X	

If bot X puts X in the bottom right position for example, it results in the following situation:

	O	
O	X	X

**Your Ultimate path for Becoming a  
DATA Scientist!**

Download this learning path to start your data science journey.

Bot O would be rejoicing (Yes! They are programmed to show emotions) as it can win the match with just one move. Now, we need to teach X not to do this again. So we give a negative reward or punishment to reinforce the correct behaviour in the next trial. We say that this action in the given state would not correspond to a negative reward and should not be considered as an optimal action in this situation.

Similarly, a positive reward would be conferred to X if it stops O from winning in the next move:

		X
	O	
O	X	

Email Id

Contact Number

 Download Resource

# Introduction to Markov Decision Process

Now that we understand the basic terminology, let's talk about formalising this whole process using a concept called a Markov Decision Process or MDP.

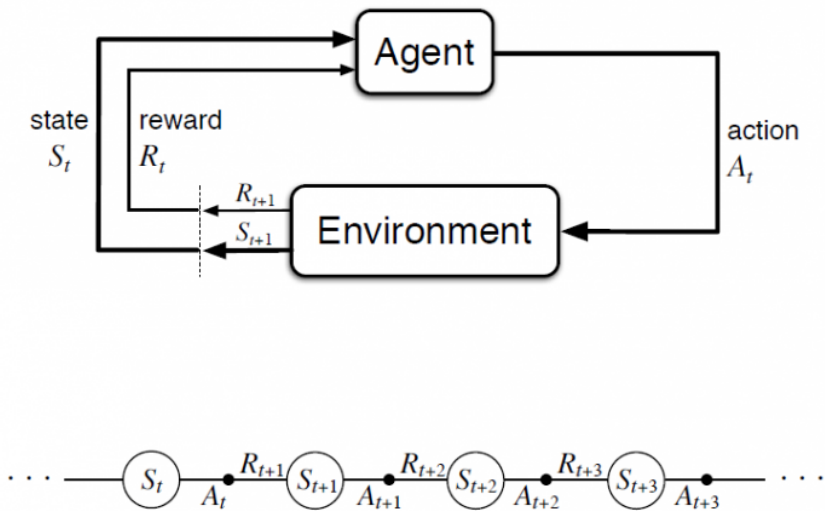
A Markov Decision Process (MDP) model contains:

- A set of possible world states  $S$
- A set of possible actions  $A$
- A real valued reward function  $R(s,a)$
- A description  $T$  of each action's effects in each state

Now, let us understand the markov or 'memoryless' property.

Any random process in which the probability of being in a given state depends only on the previous state, is a markov process.

In other words, in the markov decision process setup, the environment's response at time  $t+1$  depends only on the state and action representations at time  $t$ , and is independent of whatever happened in the past.



- $S_t$ : State of the agent at time  $t$
- $A_t$ : Action taken by agent at time  $t$
- $R_t$ : Reward obtained at time  $t$

**Your Ultimate path for Becoming a DATA Scientist!**

Download this learning path to start your data science journey.

The above diagram clearly illustrates the iteration at each time step wherein the agent receives a reward  $R_{t+1}$  and ends up in state  $S_{t+1}$  based on its action  $A_t$  at a particular state  $S_t$ . The overall goal for the agent is to maximise the cumulative reward it receives in the long run. Total reward at any time instant  $t$  is given by:

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

where  $T$  is the final time step of the episode. In the above equation, we see that all future rewards have equal weight which might not be desirable. That's where an additional concept of discounting comes into the picture. Basically, we define  $\gamma$  as a discounting factor and each reward after the immediate reward is discounted by this factor as follows:

Download Resource

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

For discount factor  $\gamma < 1$ , the rewards further in the future are getting diminished. This can be understood as a tuning parameter which can be changed based on how much one wants to consider the long term ( $\gamma$  close to 1) or short term ( $\gamma$  close to 0).

## State Value Function: How good it is to be in a given state?

Can we use the reward function defined at each time step to define how good it is, to be in a given state for a given policy? **The value function denoted as  $v(s)$  under a policy  $\pi$  represents how good a state is for an agent to be in.** In other words, what is the average reward that the agent will get starting from the current state under policy  $\pi$ ?

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right], \text{ for all } s \in \mathcal{S}$$

$\mathbb{E}$  in the above equation represents the expected reward at each state if the agent follows policy  $\pi$  and  $\mathcal{S}$  represents the set of all possible states.

Policy, as discussed earlier, is the mapping of probabilities of taking each possible action at each state ( $\pi(a/s)$ ). The policy might also be deterministic when it tells you exactly what to do at each state and does not give probabilities.

Now, it's only intuitive that 'the optimum policy' can be reached if the value function is maximised for each state. This optimal policy is then given by:

$$\pi^* = \arg \max_{\pi} V^{\pi}(s) \quad \forall s \in \mathcal{S}$$

## State-Action Value Function: How good an action is at a particular state?

The above value function only characterizes a state. Can we also know how good an action is at a particular state? A state-action value function, which is also called the q-value, does exactly that. We define the value of action  $a$ , in state  $s$ , under a policy  $\pi$ , as:

$$q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a\right]$$

This is the expected return the agent will get if it takes action  $A_t$  at time  $t$ , given state  $S_t$ , and thereafter follows policy  $\pi$ .

**Your Ultimate path for Becoming a DATA Scientist!**

Download this learning path to start your data science journey.




Download Resource

**Bellman Expectation Equation: The value information from successor states is being transferred back to the current state**



Bellman was an applied mathematician who derived equations that help to solve an Markov Decision Process.

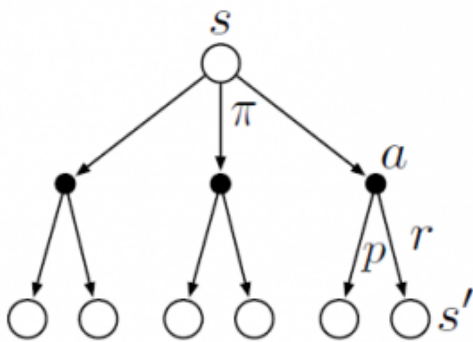
Let's go back to the state value function  $v$  and state-action value function  $q$ . Unroll the value function equation to get:

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t=s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t=s] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) \left[ r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1}=s'] \right] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) \left[ r + \gamma v_{\pi}(s') \right], \quad \text{for all } s \in \mathcal{S}, \end{aligned}$$

In this equation, we have the value function for a given policy  $\pi$  represented in terms of the value function of the next state.

Choose an action  $a$ , with probability  $\pi(a/s)$  at the state  $s$ , which leads to state  $s'$  with prob  $p(s'/s,a)$ . This gives a reward  $[r + \gamma \cdot v_{\pi}(s)]$  as given in the square bracket above.

This is called the Bellman Expectation Equation. The value information from successor states is being transferred back to the current state, and this can be represented efficiently by something called a backup diagram as shown below.



Backup diagram for  $v_{\pi}$

Your Ultimate path for Becoming a  
DATA Scientist!

The Bellman expectation equation averages over all the possibilities, weighting each by its probability of occurring. It states that the value of the start state must equal the (discounted) value of the expected next state, plus the reward expected along the way.

We have  $n$  (number of states) linear equations with unique solution to solve for each state  $s$ .

Bellman Optimality Equation: Find the optimal policy

The goal here is to find the optimal policy, which when followed by the agent gets the maximum cumulative reward. In other words, find a policy  $\pi$ , such that for no other  $\pi$  can the agent get a better expected return. We want to find a policy which achieves maximum value for each state.

Name

Email Id

Contact Number

Download Resource

$$\pi^* = \arg \max_{\pi} V^{\pi}(s) \quad \forall s \in \mathbb{S}$$

Note that we might not get a unique policy, as under any situation there can be 2 or more paths that have the same return and are still optimal.

$$v_*(s) \doteq \max_{\pi} v_{\pi}(s)$$

$$q_*(s, a) \doteq \max_{\pi} q_{\pi}(s, a)$$

Optimal value function can be obtained by finding the action  $a$  which will lead to the maximum of  $q^*$ . This is called the bellman optimality equation for  $v^*$ .

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

Intuitively, the Bellman optimality equation says that the value of each state under an optimal policy must be the return the agent gets when it follows the best action as given by the optimal policy. For optimal policy  $\pi^*$ , the optimal value function is given by:

$$\begin{aligned} v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\ &= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

Given a value function  $q^*$ , we can recover an optimum policy as follows:

$$\begin{aligned} \pi'(s) &\doteq \arg \max_a q_{\pi}(s, a) \\ &= \arg \max_a \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

**Your Ultimate path for Becoming a  
DATA Scientist!**

The value function for optimal policy can be solved through a non-linear system of equations. We can solve these efficiently using iterative methods that fall under the umbrella of dynamic programming.

## Dynamic Programming

Dynamic programming algorithms solve a category of problems called planning problems. Herein given the complete model and specifications of the environment (MDP), we can successfully find an optimal policy for the agent to follow. It contains two main steps:

1. Break the problem into subproblems and solve it
2. Solutions to subproblems are cached or stored for reuse to find overall optimal solution to the problem at hand



To solve a given MDP, the solution must have the components to:

- 1. Find out how good an arbitrary policy is
- 2. Find out the optimal policy for the given MDP


**Policy Evaluation: Find out how good a policy is?**


Policy evaluation answers the question of how good a policy is. Given an MDP and an arbitrary policy  $\pi$ , we will compute the state-value function. This is called policy evaluation in the DP literature. The idea is to turn bellman expectation equation discussed earlier to an update.

$$v_{k+1}(s) \doteq \mathbb{E}_{\pi}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s]$$

To produce each successive approximation  $v_{k+1}$  from  $v_k$ , iterative policy evaluation applies the same operation to each state  $s$ . It replaces the old value of  $s$  with a new value obtained from the old values of the successor states of  $s$ , and the expected immediate rewards, along all the one-step transitions possible under the policy being evaluated, until it converges to the true value function of a given policy  $\pi$ .

Let us understand policy evaluation using the very popular example of Gridworld.


1	2	3	4
5	6 	7	8
9	10	11	12
13	14	15	16

actions 

Reward is -1 for  
all transition

A bot is required to traverse a grid of 4x4 dimensions to reach its goal. Each step is associated with a reward of -1. There are 2 terminal states here: 1 and 16 and 14 non-terminal states given by 2,3,4,5,6,7,8,9,10,11,12,13,15. We will start with initialising  $v_0$  for the random policy to all 0s.

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

 Download Resource

Your Ultimate path for Becoming a  
DATA Scientist!

States given by 2,3,4,5,6,7,8,9,10,11,12,13,15  
Consider a random policy for which, at every state, the probability of every action (up, down, left, right) is equal to 0.25. We will start with

This is definitely not very useful. Let's calculate  $v_2$  for all the states of 6:

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

$$\begin{aligned}
 v_1(6) &= \sum_{a \in \{u,d,l,r\}} \pi(a|6) \sum_{s',r} p(s',r|6,a) [r + \gamma v_0(s')] \\
 &= \sum_{a \in \{u,d,l,r\}} \underbrace{\pi(a|6)}_{= 0.25 \forall a} \sum_{s'} p(s'|6,a) \underbrace{[r + \gamma v_0(s')]}_{\substack{= -1 \\ = 0 \forall s'}} \\
 &= 0.25 * \{-p(2|6,u) - p(10|6,d) - p(5|6,l) - p(7|6,r)\} \\
 &= 0.25 * \{-1 - 1 - 1 - 1\} \\
 &= -1 \\
 &\Rightarrow v_1(6) = -1
 \end{aligned}$$

Similarly, for all non-terminal states,  $v_1(s) = -1$ .

For terminal states  $p(s'/s,a) = 0$  and hence  $v_k(1) = v_k(16) = 0$  for all  $k$ . So  $v_1$  for the random policy is given by:

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

Now, for  $v_2(s)$  we are assuming  $\gamma$  or the discounting factor to be 1:

## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.

$$\begin{aligned}
 v_2(6) &= \sum_{a \in \{u,d,l,r\}} \underbrace{\pi(a|6)}_{= 0.25 \forall a} \sum_{s'} p(s'|6, a) \underbrace{[r + \gamma v_1(s')]}_{= -1} \\
 &= 0.25 * \{p(2|6, u)[-1 - \gamma] + p(10|6, d)[-1 - \gamma] \\
 &\quad + p(5|6, l)[-1 - \gamma] + p(7|6, r)[-1 - \gamma]\} \\
 &\stackrel{\gamma=1}{=} 0.25 * \{-2 - 2 - 2 - 2\} \\
 &= -2
 \end{aligned}$$

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

As you can see, all the states marked in red in the above diagram are identical to 6 for the purpose of calculating the value function. Hence, for all these states,  $v_2(s) = -2$ .

For all the remaining states, i.e., 2, 5, 12 and 15,  $v_2$  can be calculated as follows:


$$\begin{aligned}
 v_2(2) &= \sum_{a \in \{u,d,l,r\}} \underbrace{\pi(a|2)}_{= 0.25 \forall a} \sum_{s'} p(s'|2, a) \underbrace{[r + \gamma v_1(s')]}_{= -1} \\
 &= 0.25 * \{p(2|2, u)[-1 - \gamma] + p(6|2, d)[-1 - \gamma] \\
 &\quad + p(1|2, l)[-1 - \gamma * 0] + p(3|2, r)[-1 - \gamma]\} \\
 &\stackrel{\gamma=1}{=} 0.25 * \{-2 - 2 - 1 - 2\} \\
 &= -1.75 \\
 &\Rightarrow v_2(2) = -1.75
 \end{aligned}$$

**Your Ultimate path for Becoming a DATA Scientist!**

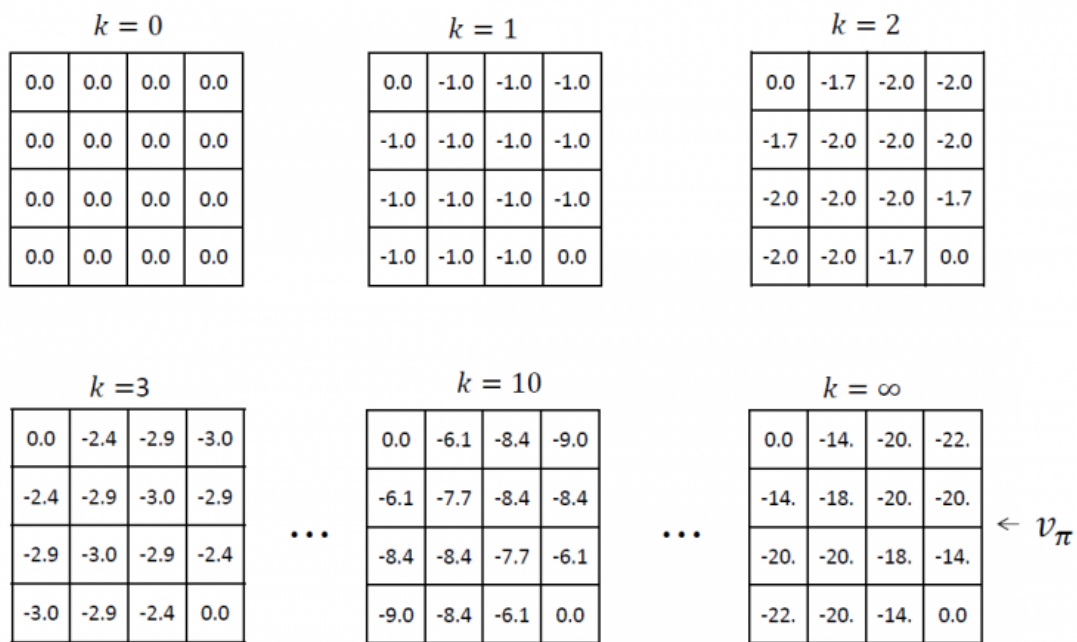
Download this learning path to start your data science journey.

$\Rightarrow v_2$  for the random policy:

0.0	Name	-2.0	-2.0	
-1.7	Email Id	-2.0	-2.0	
-2.0	Contact Number	-2.0	-1.7	
-2.0		-2.0	-1.7	0.0

 Download Resource

If we repeat this step several times, we get  $v_\pi$ .



## Policy Improvement: Improve an arbitrary policy

Using policy evaluation we have determined the value function  $v$  for an arbitrary policy  $\pi$ . We know how good our current policy is. Now for some state  $s$ , we want to understand what is the impact of taking an action  $a$  that does not pertain to policy  $\pi$ . Let's say we select  $a$  in  $s$ , and after that we follow the original policy  $\pi$ . The value of this way of behaving is represented as:

$$\begin{aligned}
 q_\pi(s, a) &\doteq \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\
 &= \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_\pi(s')].
 \end{aligned}$$

If this happens to be greater than the value function  $v_\pi(s)$ , it implies that the new policy  $\pi'$  would be better to take. We do this iteratively for all states to find the best policy. Note that in this case, the agent would be following a greedy policy in the sense that it is looking only one step ahead.

## Your Ultimate path for Becoming a DATA Scientist!

Let's get back to our example of gridworld. Using  $v_\pi$ , the value function obtained for a random policy  $\pi$ , we can improve upon  $\pi$  by following the path of highest value (as shown in the figure below). We start with an arbitrary policy, and for each state one step look-ahead is done to find the action leading to the state with the highest value. This is done successively for each state.

As shown below for state 2, the optimal action is left which leads to the terminal state having a value . This is the highest among all the next states (0,-18,-20). This is repeated for all states to find the new policy.

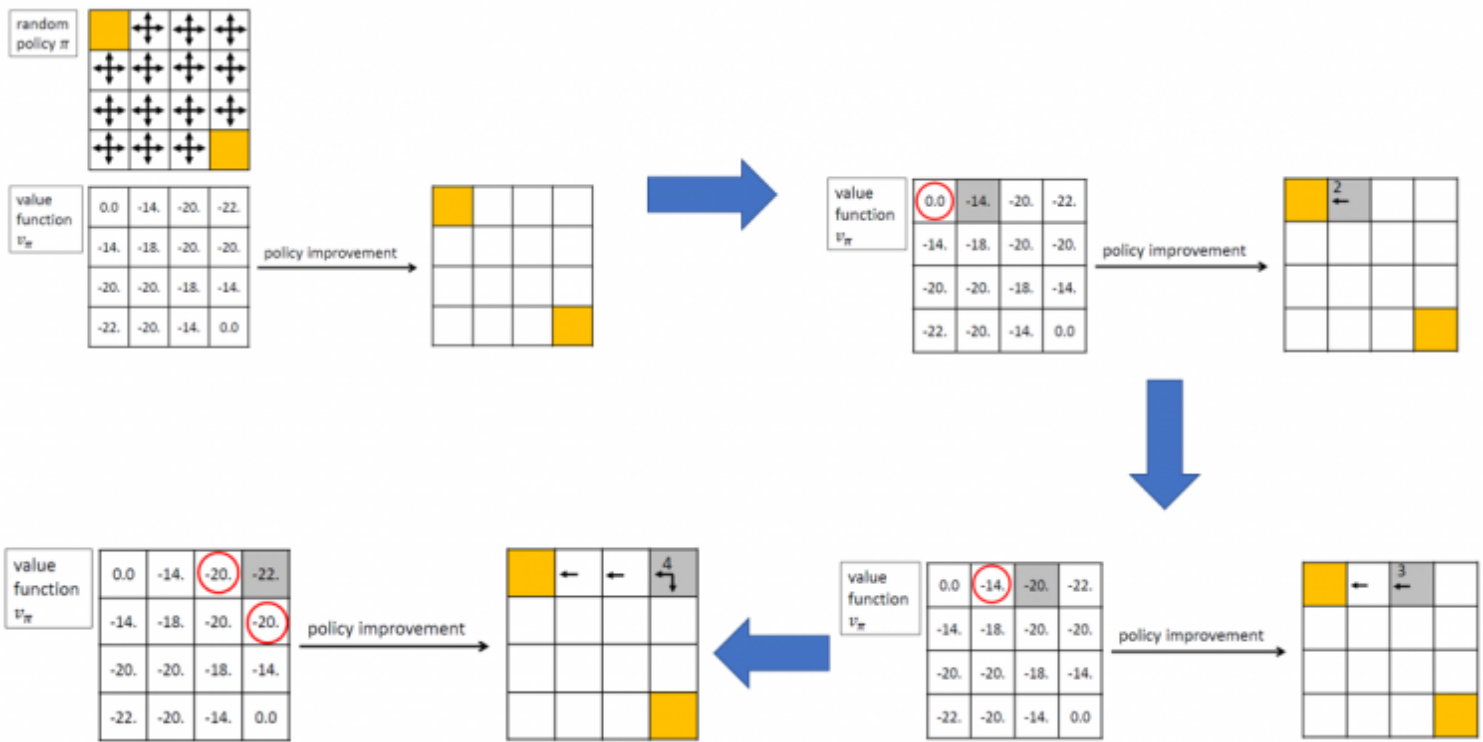
Name

Email Id

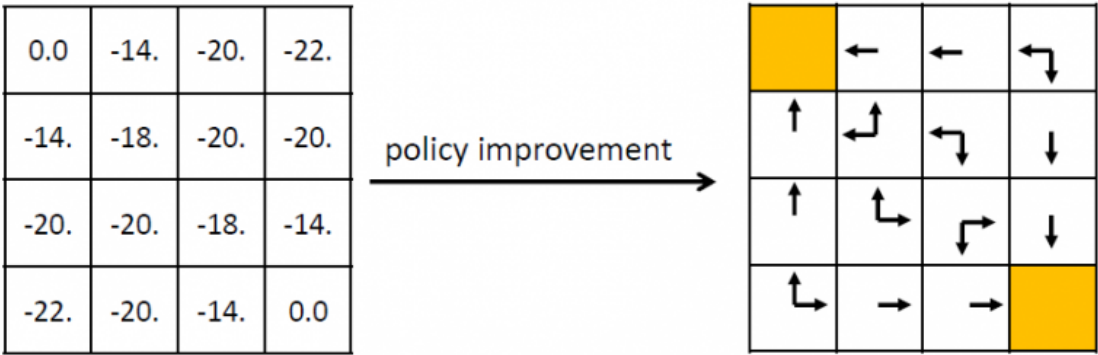
Contact Number



Download Resource



Overall, after the policy improvement step using  $v_\pi$ , we get the new policy  $\pi'$ :



Looking at the new policy, it is clear that it's much better than the random policy. However, we should calculate  $v_{\pi'}$  using the policy evaluation technique we discussed earlier to verify this claim.

**Your Ultimate path for Becoming a DATA Scientist!**

Download this learning path to start your data science journey.

### Policy Iteration: Policy Evaluation + Policy Improvement

Once the policy has been improved using  $v_\pi$  to yield a better policy  $\pi'$ , we can then compute  $v_{\pi'}$  to improve it further to  $\pi''$ . Repeated iterations are done to converge approximately to the true value function for a given policy  $\pi$  (policy evaluation). Improving the policy as described in the policy improvement section is called policy iteration.

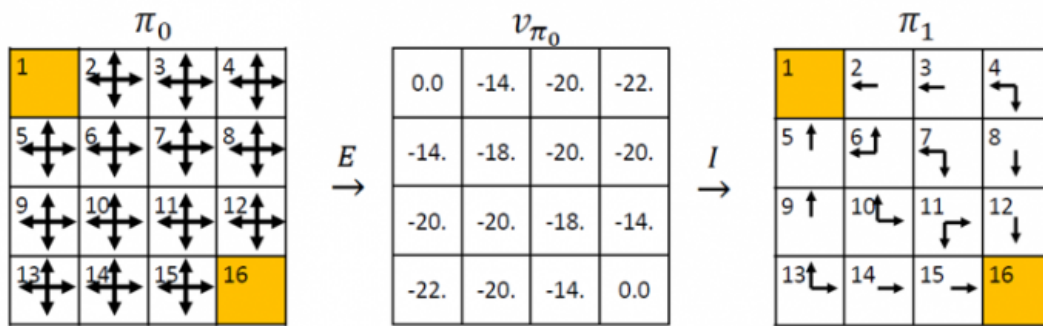
Name

Email Id

Contact Number

Download Resource

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \rightarrow \dots \rightarrow \pi_* \xrightarrow{E} v_*$$



In this way, the new policy is sure to be an improvement over the previous one and given enough iterations, it will return the optimal policy. This sounds amazing but there is a drawback – each iteration in policy iteration itself includes another iteration of policy evaluation that may require multiple sweeps through all the states. Value iteration technique discussed in the next section provides a possible solution to this.

## Value Iteration

We saw in the gridworld example that at around  $k = 10$ , we were already in a position to find the optimal policy. So, instead of waiting for the policy evaluation step to converge exactly to the value function  $v_{\pi}$ , we could stop earlier.

We can also get the optimal policy with just 1 step of policy evaluation followed by updating the value function repeatedly (but this time with the updates derived from bellman optimality equation). Let's see how this is done as a simple backup operation:

$$v_{k+1}(s) \doteq \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a]$$

This is identical to the bellman update in policy evaluation, with the difference being that we are taking the maximum over all actions. Once the updates are small enough, we can take the value function obtained as final and estimate the optimal policy corresponding to that.

Some important points related to DP:

1. DP can only be used if the model of the environment is known.
2. Has a very high computational expense, i.e., it does not scale well as the number of states increase to a large number. An alternative called asynchronous dynamic programming helps to resolve this issue to some extent.

## DP in action: Finding optimal policy for Frozen Lake environment using Python

It is of utmost importance to first have a defined environment in order to test any kind of policy for solving an MDP efficiently. Thankfully, OpenAI, a non profit research organization provides a large number of environments to test and play with various reinforcement learning algorithms. To illustrate dynamic programming here, we will use it to

**Your Ultimate path for Becoming a  
DATA Scientist!**

Download this learning path to start your data  
science journey.

Name

Email Id

Contact Number



Download Resource



navigate the Frozen Lake environment.

## Frozen Lake Environment

The agent controls the movement of a character in a grid world. Some tiles of the grid are walkable, and others lead to the agent falling into the water. Additionally, the movement direction of the agent is uncertain and only partially depends on the chosen direction. The agent is rewarded for finding a walkable path to a goal tile.

The surface is described using a grid like the following:

S	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

(S: starting point, safe), (F: frozen surface, safe), (H: hole, fall to your doom), (G: goal)

The idea is to reach the goal from the starting point by walking only on frozen surface and avoiding all the holes. Installation details and documentation is available at this [link \(https://gym.openai.com/docs/\)](https://gym.openai.com/docs/).

Once gym library is installed, you can just open a jupyter notebook to get started.

```
import gym
import numpy as np
env = gym.make('FrozenLake-v0')
```

Now, the env variable contains all the information regarding the frozen lake environment. Before we move on, we need to understand what an episode is. An episode represents a trial by the agent in its pursuit to reach the goal. An episode ends once the agent reaches a terminal state which in this case is either a hole or the goal.

## Policy Iteration in python

Description of parameters for policy iteration function

**policy:** 2D array of a size  $n(S) \times n(A)$ , each cell represents a probability of taking action  $a$  in state  $s$ .

**environment:** Initialized OpenAI gym environment object

**discount\_factor:** MDP discount factor

**theta:** A threshold of a value function change. Once the update to value function is below this number

**max\_iterations:** Maximum number of iterations to avoid letting the program run indefinitely

## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.

Download Resource

This function will return a vector of size  $nS$ , which represent a value function for each state.

Let's start with the policy evaluation step. The objective is to converge to the true value function for a given policy  $\pi$ . We will define a function that returns the required value function.

```
def policy_evaluation(policy, environment, discount_factor=1.0, theta=1e-9, max_iterations=1e9):
    # Number of evaluation iterations
    evaluation_iterations = 1
    # Initialize a value function for each state as zero
    V = np.zeros(environment.nS)
    # Repeat until change in value is below the threshold
    for i in range(int(max_iterations)):
        # Initialize a change of value function as zero
        delta = 0
        # Iterate though each state
        for state in range(environment.nS):
            # Initial a new value of current state
            v = 0
            # Try all possible actions which can be taken from this state
            for action, action_probability in enumerate(policy[state]):
                # Check how good next state will be
                for state_probability, next_state, reward, terminated in environment.P[state][action]:
                    # Calculate the expected value
                    v += action_probability * state_probability * (reward + discount_factor * V[next_state])

            # Calculate the absolute change of value function
            delta = max(delta, np.abs(V[state] - v))
            # Update value function
            V[state] = v
        evaluation_iterations += 1

    # Terminate if value change is insignificant
    if delta < theta:
        print(f'Policy evaluated in {evaluation_iterations} iterations.')
        return V
```

## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.

Name

Email Id

Contact Number



Download Resource

Now coming to the **policy improvement** part of the policy iteration algorithm. We need a helper function that does one step lookahead to calculate the state-value function. This will return an array of length  $nA$  containing expected value of each action

```
def one_step_lookahead(environment, state, V, discount_factor):
    action_values = np.zeros(environment.nA)
    for action in range(environment.nA):
        for probability, next_state, reward, terminated in environment.P[state][action]:
            action_values[action] += probability * (reward + discount_factor * V[next_state])
    return action_values
```

Now, the overall policy iteration would be as described below. This will return a tuple (policy,V) which is the optimal policy matrix and value function for each state.

```
def policy_iteration(environment, discount_factor=1.0, max_iterations=1e9):
    # Start with a random policy
    #num states x num actions / num actions
    policy = np.ones([environment.nS, environment.nA]) / environment.nA
    # Initialize counter of evaluated policies
    evaluated_policies = 1
    # Repeat until convergence or critical number of iterations reached
    for i in range(int(max_iterations)):
        stable_policy = True
        # Evaluate current policy
        V = policy_evaluation(policy, environment, discount_factor=discount_factor)
        # Go through each state and try to improve actions that were taken (policy Improvement)
        for state in range(environment.nS):
            # Choose the best action in a current state under current policy
            current_action = np.argmax(policy[state])
            # Look one step ahead and evaluate if current action is optimal
            # We will try every possible action in a current state
            action_value = one_step_lookahead(environment, state, V, discount_factor)
            # Select a better action
            best_action = np.argmax(action_value)
            # If action didn't change
            if current_action != best_action:
                stable_policy = False
            # Greedy policy update
            policy[state] = np.eye(environment.nA)[best_action]
        evaluated_policies += 1
    # If the algorithm converged and policy is not changing anymore, then return final policy and
    value function
    if stable_policy:
        print(f'Evaluated {evaluated_policies} policies.')
        return policy, V
```

**Your Ultimate path for Becoming a  
DATA Scientist!**

Download this learning path to start your data  
science journey.





Download Resource

## Value Iteration in python

The parameters are defined in the same manner for value iteration. The value iteration algorithm can be similarly coded:

```
def value_iteration(environment, discount_factor=1.0, theta=1e-9, max_iterations=1e9):
    # Initialize state-value function with zeros for each environment state
    V = np.zeros(environment.nS)
    for i in range(int(max_iterations)):
        # Early stopping condition
        delta = 0
        # Update each state
        for state in range(environment.nS):
            # Do a one-step lookahead to calculate state-action values
            action_value = one_step_lookahead(environment, state, V, discount_factor)
            # Select best action to perform based on the highest state-action value
            best_action_value = np.max(action_value)
            # Calculate change in value
            delta = max(delta, np.abs(V[state] - best_action_value))
            # Update the value function for current state
            V[state] = best_action_value
        # Check if we can stop
        if delta < theta:
            print(f'Value-iteration converged at iteration#{i}.')
            break

    # Create a deterministic policy using the optimal value function
    policy = np.zeros([environment.nS, environment.nA])
    for state in range(environment.nS):
        # One step lookahead to find the best action for this state
        action_value = one_step_lookahead(environment, state, V, discount_factor)
        # Select best action based on the highest state-action value
        best_action = np.argmax(action_value)
        # Update the policy to perform a better action at a current state
        policy[state, best_action] = 1.0

    return policy, V
```

**Your Ultimate path for Becoming a  
DATA Scientist!**

Download this learning path to start your data  
science journey.



Download Resource

Finally, let's compare both methods to look at which of them works better in a practical setting. To do this, we will try to learn the optimal policy for the frozen lake environment using both techniques described above. Later, we will check which technique performed better based on the average return after 10,000 episodes.

```
def play_episodes(environment, n_episodes, policy):
    wins = 0
    total_reward = 0
    for episode in range(n_episodes):
        terminated = False
        state = environment.reset()
        while not terminated:
            # Select best action to perform in a current state
            action = np.argmax(policy[state])
            # Perform an action and observe how environment acted in response
            next_state, reward, terminated, info = environment.step(action)
            # Summarize total reward
            total_reward += reward
            # Update current state
            state = next_state
            # Calculate number of wins over episodes
            if terminated and reward == 1.0:
                wins += 1
        average_reward = total_reward / n_episodes
    return wins, total_reward, average_reward

# Number of episodes to play
n_episodes = 10000

# Functions to find best policy
solvers = [('Policy Iteration', policy_iteration),
            ('Value Iteration', value_iteration)]

for iteration_name, iteration_func in solvers:
    # Load a Frozen Lake environment
    environment = gym.make('FrozenLake-v0')
    # Search for an optimal policy using policy iteration
    policy, V = iteration_func(environment.env)
    # Apply best policy to the real environment
    wins, total_reward, average_reward = play_episodes(environment, n_episodes, policy)
    print(f'{iteration_name} :: number of wins over {n_episodes} episodes = {wins}')
    print(f'{iteration_name} :: average reward over {n_episodes} episodes = {average_reward} \n\n')
    Name
```

Policy evaluated in 66 iterations.

Evaluated 2 policies.

Policy Iteration :: number of wins over 10000 episodes = 7287

Policy Iteration :: average reward over 10000 episodes = 0.7287

Value-iteration converged at iteration#523.

Value Iteration :: number of wins over 10000 episodes = 7397

Value Iteration :: average reward over 10000 episodes = 0.7397

## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey

Email Id

Contact Number



Download Resource

We observe that value iteration has a better average reward and higher number of wins when it is run for 10,000 episodes.

## End Notes

In this article, we became familiar with model based planning using dynamic programming, which given all specifications of an environment, can find the best policy to take. I want to particularly mention the brilliant book on RL by Sutton and Barto which is a bible for this technique and encourage people to refer it. More importantly, you have taken the first step towards mastering reinforcement learning. Stay tuned for more articles covering different algorithms within this exciting domain.

You can also read this article on Analytics Vidhya's Android APP



[https://play.google.com/store/apps/details?id=com.analyticsvidhya.android&utm\\_source=blog\\_article&utm\\_campaign=blog&pcampaignid=MKT-Other-global-all-co-prtnr-py-PartBadge-Mar2515-1](https://play.google.com/store/apps/details?id=com.analyticsvidhya.android&utm_source=blog_article&utm_campaign=blog&pcampaignid=MKT-Other-global-all-co-prtnr-py-PartBadge-Mar2515-1)


### Share this:

 (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?share=linkedin&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?share=facebook&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?share=twitter&nb=1>)

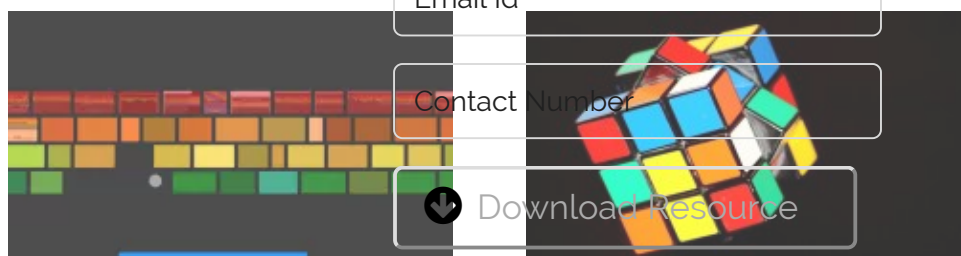
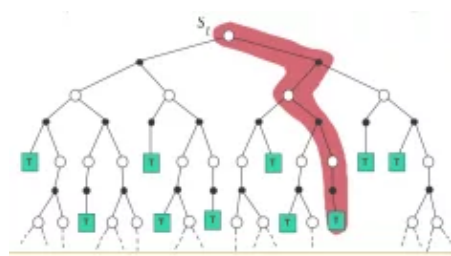
 (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?share=pocket&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?share=reddit&nb=1>)

### Like this:

Loading...

## Related Articles



## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.



<a href="https://www.analyticsvidhya.com/blog/2018/11/reinforcement-learning-introduction-monte-carlo-learning-openai-gym/">https://www.analyticsvidhya.com/blog/2018/11/reinforcement-learning-introduction-monte-carlo-learning-openai-gym/</a> Reinforcement Learning: Introduction to Monte Carlo Learning using the OpenAI Gym Toolkit <a href="https://www.analyticsvidhya.com/blog/2018/11/reinforcement-learning-introduction-monte-carlo-learning-openai-gym/">https://www.analyticsvidhya.com/blog/2018/11/reinforcement-learning-introduction-monte-carlo-learning-openai-gym/</a> November 19, 2018 In "Python"	<a href="https://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/">https://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/</a> A Hands-On Introduction to Deep Q-Learning using OpenAI Gym in Python <a href="https://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/">https://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/</a> April 18, 2019 In "Python"	<a href="https://www.analyticsvidhya.com/blog/2017/01/introduction-to-reinforcement-learning-implementation/">https://www.analyticsvidhya.com/blog/2017/01/introduction-to-reinforcement-learning-implementation/</a> Simple Beginner's guide to Reinforcement Learning & its implementation <a href="https://www.analyticsvidhya.com/blog/2017/01/introduction-to-reinforcement-learning-implementation/">https://www.analyticsvidhya.com/blog/2017/01/introduction-to-reinforcement-learning-implementation/</a> January 19, 2017 In "Machine Learning"
---	--	---

TAGS : [PYTHON \(https://www.analyticsvidhya.com/blog/tag/python/\)](https://www.analyticsvidhya.com/blog/tag/python/), [REINFORCEMENT LEARNING \(https://www.analyticsvidhya.com/blog/tag/reinforcement-learning/\)](https://www.analyticsvidhya.com/blog/tag/reinforcement-learning/)

NEXT ARTICLE

### Performing Speech and Object Recognition using just One Model with MIT's ML System

<https://www.analyticsvidhya.com/blog/2018/09/speech-object-recognition-one-model-mit-ml/>

...

PREVIOUS ARTICLE

### DataHack Radio #10: The Role of Computer Science in the Data Science World with Dr. Jeannette M. Wing

<https://www.analyticsvidhya.com/blog/2018/09/datahack-radio-data-science-podcast-jeannette-wing/>



### Ankit Choudhary [\(https://www.analyticsvidhya.com/blog/author/ankit2106/\)](https://www.analyticsvidhya.com/blog/author/ankit2106/) **Your Ultimate path for Becoming a DATA Scientist!**

IIT Bombay Graduate with a Masters and Bachelors in Electrical Engineering. I have previously worked as a lead decision scientist for Indian National Congress deploying statistical models (Segmentation, K-Nearest Neighbours) to help party leadership/Team make data-driven decisions. My interest lies in putting data in heart of business for data-driven decision making.

**in** <https://www.linkedin.com/in/ankit-choudhary-b9360826/>

Name

Email Id

Contact Number

This article is quite old and you might not get a prompt response from the author. We request you to post this comment on Analytics Vidhya's **Discussion portal** <https://discuss.analyticsvidhya.com/> to get your queries resolved



Download Resource

## 4 COMMENTS



**RAMESH MATHIKUMAR**

[Reply](#)

September 18, 2018 at 11:11 pm (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/#comment-155008>)

Good Article Ankit.



**BHUMIKA**

[Reply](#)

September 24, 2018 at 11:55 am (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/#comment-155059>)

Hello. How do we derive the Bellman expectation equation?



**ANKIT CHOUDHARY**

[Reply](#)

September 24, 2018 at 12:20 pm (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/#comment-155060>)

You can refer to this stack overflow query: <https://stats.stackexchange.com/questions/243384/deriving-bellmans-equation-in-reinforcement-learning> (<https://stats.stackexchange.com/questions/243384/deriving-bellmans-equation-in-reinforcement-learning>) for the derivation.



**VIJIT**

[Reply](#)

April 3, 2019 at 1:18 pm (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/#comment-157693>)

Excellent article on Dynamic Programming. Explained the concepts in a very easy way.

## JOIN THE NEXTGEN DATA SCIENCE ECOSYSTEM

Get access to free courses on Analytics Vidhya

Get free downloadable resource from Analytics Vidhya

Save your articles

Participate in hackathons and win prizes

([https://id.analyticsvidhya.com/accounts/login/?next=https://www.analyticsvidhya.com/blog/?utm\\_source=blog-subscribe&utm\\_medium=web](https://id.analyticsvidhya.com/accounts/login/?next=https://www.analyticsvidhya.com/blog/?utm_source=blog-subscribe&utm_medium=web))

## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.

Join Now

 Download Resource



([https://software.seek.intel.com/DataCenter\\_to\\_Edge\\_REG?](https://software.seek.intel.com/DataCenter_to_Edge_REG?registration_source=AnalyticsVidhya-APJ)

registration\_source=AnalyticsVidhya-APJ)

## POPULAR POSTS

24 Ultimate Data Science Projects To Boost Your Knowledge and Skills (& can be accessed freely)

(<https://www.analyticsvidhya.com/blog/2018/05/24-ultimate-data-science-projects-to-boost-your-knowledge-and-skills/>)

Essentials of Machine Learning Algorithms (with Python and R Codes)

(<https://www.analyticsvidhya.com/blog/2017/09/common-machine-learning-algorithms/>)

7 Types of Regression Techniques you should know!

(<https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/>)

A Complete Tutorial to Learn Data Science with Python from Scratch

(<https://www.analyticsvidhya.com/blog/2016/01/complete-tutorial-learn-data-science-python-scratch-2/>)

Understanding Support Vector Machine algorithm from examples (along with code)

(<https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/>)

Stock Prices Prediction Using Machine Learning and Deep Learning Techniques (with Python codes)

(<https://www.analyticsvidhya.com/blog/2018/10/predicting-stock-prices-using-machine-learning-and-deep-learning-techniques-python/>)

Introduction to k-Nearest Neighbors: Simplified (with implementation in Python)

(<https://www.analyticsvidhya.com/blog/2018/03/introduction-k-neighbours-algorithm-clustering/>)

A Simple Introduction to ANOVA (with applications in Excel)

(<https://www.analyticsvidhya.com/blog/2018/01/anova-analysis-of-variance/>)

**Your Ultimate path for Becoming a  
DATA Scientist!**

Download this learning path to start your data science journey

## RECENT POSTS

Top 7 Machine Learning Github Repositories for Data Scientists (<https://www.analyticsvidhya.com/blog/2019/06/top-7-machine-learning-github-repositories-data-scientists/>)

JUNE 6, 2019

The AI Comic: Z.A.I.N – Issue #1: Automating Attendance using Computer Vision (<https://www.analyticsvidhya.com/blog/2019/06/ai-comic-zain-issue-1-automating-computer-vision/>)

JUNE 3, 2019

DataHack Radio #23: Ines Montani and Matthew Honnibal – The Brains behind spaCy (<https://www.analyticsvidhya.com/blog/2019/06/datahack-radio-ines-montani-matthew-honnibal-brains-behind-spacy/>)

JUNE 3, 2019

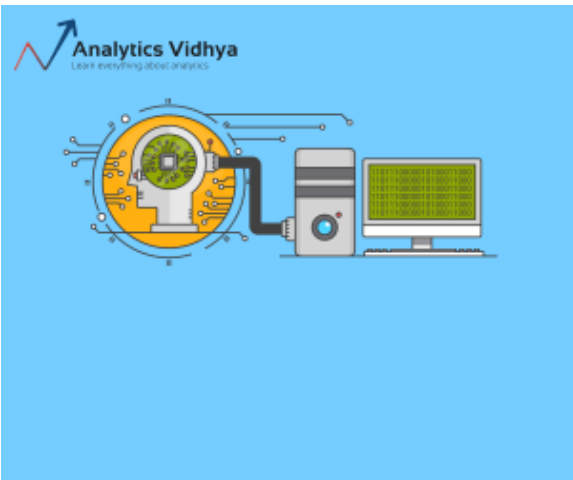
Exclusive Interview with Sonny Laskar – Kaggle Master and Analytics Vidhya Hackathon Expert (<https://www.analyticsvidhya.com/blog/2019/05/exclusive-interview-sonny-laskar-kaggle-master-analytics-vidhya-hackathon-expert/>)

MAY 30, 2019



([https://datahack.analyticsvidhya.com/contest/intel-ai-enterprise-](https://datahack.analyticsvidhya.com/contest/intel-ai-enterprise-meetup-mumbai-invite-only/?utm_source=Sticky_banner1&utm_medium=display&utm_campaign=IntelMum)

[meetup-mumbai-invite-only/?utm\\_source=Sticky\\_banner1&utm\\_medium=display&utm\\_campaign=IntelMum](https://datahack.analyticsvidhya.com/contest/intel-ai-enterprise-meetup-mumbai-invite-only/?utm_source=Sticky_banner1&utm_medium=display&utm_campaign=IntelMum))



([https://courses.analyticsvidhya.com/courses/applied-machine-](https://courses.analyticsvidhya.com/courses/applied-machine-learning-beginner-to-professional?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=Applied_ML)  
[learning-beginner-to-professional?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=Applied\\_ML](https://courses.analyticsvidhya.com/courses/applied-machine-learning-beginner-to-professional?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=Applied_ML))

## Your Ultimate path for Becoming a DATA Scientist!

Analytics Vidhya has a learning path to start your data science journey.

 Download Resource

## ANALYTICS VIDHYA

About Us

(http://www.analyticsvidhya.com/about-me/)

Our Team

(https://www.analyticsvidhya.com/about-me/team/)

Career

(https://www.analyticsvidhya.com/career-analytcs-vidhya/)

Contact Us

(https://www.analyticsvidhya.com/contact/)

Write for us

(https://www.analyticsvidhya.com/about-me/write/)

## DATA SCIENTISTS

Blog

(https://www.analyticsvidhya.com/blog/)

Hackathon

(https://datahack.analyticsvidhya.com/)

Discussions

(https://discuss.analyticsvidhya.com/)

Apply Jobs

(https://www.analyticsvidhya.com/contact/)

Leaderboard

(https://datahack.analyticsvidhya.com/leaderboard/)

## COMPANIES

Post Jobs

(https://www.analyticsvidhya.com/corporate/)

Trainings

(https://trainings.analyticsvidhya.com/)

Hiring Hackathons

(https://www.analyticsvidhya.com/hiring-hackathons/)

Advertising

(https://www.analyticsvidhya.com/contact/)

Reach Us

(https://www.analyticsvidhya.com/contact/)

## JOIN OUR COMMUNITY :

f

(https://www.facebook.com/analyticsvidhya)

🐦

(https://twitter.com/analyticsvidhya)

43356

20938

(https://www.facebook.com/analyticsvidhya)

(https://twitter.com/analyticsvidhya)

Followers

(https://www.facebook.com/analyticsvidhya)

Followers

(https://twitter.com/analyticsvidhya)

Followers

(https://www.facebook.com/analyticsvidhya)

in

(https://www.linkedin.com/company/analyticsvidhya)

(https://plus.google.com/+Analyticsvidhya)

(https://plus.google.com/+Analyticsvidhya)

Followers

(https://plus.google.com/+Analyticsvidhya)

(https://in.linkedin.com/company/analyticsvidhya)

(https://in.linkedin.com/company/analyticsvidhya)

(https://plus.google.com/+Analyticsvidhya)

(https://plus.google.com/+Analyticsvidhya)

(https://plus.google.com/+Analyticsvidhya)

vidhya)

Subscribe to emailer

>

© Copyright 2013-2019 Analytics Vidhya.

Privacy Policy (https://www.analyticsvidhya.com/privacy-policy/)

Terms of Use (https://www.analyticsvidhya.com/terms/)

Refund Policy (https://www.analyticsvidhya.com/refund-policy/)

Don't have an account? Sign up (https://id.analyticsvidhya.com/)

x

-

(http://play.google.com/store/apps/details?id=com.analyticsvidhya.android)

## Your Ultimate path for Becoming a DATA Scientist!

Download this learning path to start your data science journey.

Name

Email Id

Contact Number



Download Resource