

x

-

BLOG ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/blog/?utm_source=home_blog_navbar)) ▾



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

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ▾

HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))



DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar))

G+ ([HTTPS://PLUS.GOOGLE.COM](https://plus.google.com))

in ([HTTPS://IN.LINKEDIN.COM](https://in.linkedin.com))

DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=home_blog_navbar))



LOGIN / REGISTER ([HTTPS://ID.ANALYTICSVIDHYA.COM](https://id.analyticsvidhya.com))

NEXT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2018/09/REINFORCEMENT-LEARNING-MODEL-BASED-PLANNING-DYNAMIC-PROGRAMMING/](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/))

utm_source=home_blog_navbar)

CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT](https://www.analyticsvidhya.com/contact))

HOME ([HTTPS://WWW.ANALYTICSVIDHYA.COM](https://www.analyticsvidhya.com)) BLOG ARCHIVE ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG-ARCHIVE](https://www.analyticsvidhya.com/blog-archive))

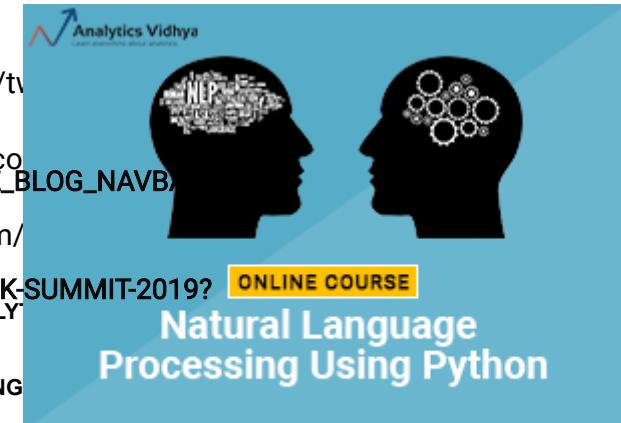
DISCUSS ([HTTPS://DISCUSS.ANALYTICSVIDHYA.COM](https://discuss.analyticsvidhya.com)) CORPORATE ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE](https://www.analyticsvidhya.com/corporate))



([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=TOPBANNER&UTM_MEDIUM=TOPBANNER](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=blog&utm_medium=topBanner))

[Home](https://www.analyticsvidhya.com/) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/](https://www.analyticsvidhya.com/)) > Python

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY PYTHON-2/](https://www.analyticsvidhya.com/blog/category/python-2/)) > Nuts & Bolts of



Learn to Solve
Text Classification Problems Using **NLP**

([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/NATURAL-LANGUAGE-PROCESSING-NLP/](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/))



Start Building Your First
Computer Vision Model Today

Reinforcement Learning: Model Based Planning using Dynamic Programming BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR) (https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/).

COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ▾

PYTHON (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/PYTHON-2/)

REINFORCEMENT LEARNING

HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL) (https://www.analyticsvidhya.com/blog/category/machine-learning/reinforcement-learning/)

DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR)

Nuts & Bolts of Reinforcement Learning:

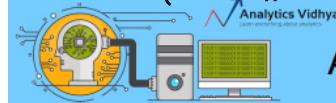
Model Based Planning using Dynamic Programming

Programming

UTM_SOURCE=HOME_BLOG_NAVBAR)

ANKIT CHAUDHARY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/AUTHOR/ANKIT21...)

CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)



Applied Machine Learning

BEGINNER TO PROFESSIONAL

JOIN NOW

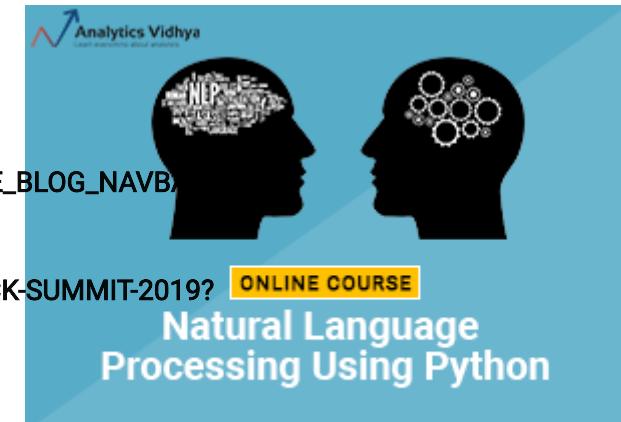
ONLINE COURSE

(https://courses.analyticsvidhya.com/courses/applied-machine-learning-beginner-to-professional/?utm_source=blog&utm_medium=bannerbelowblog).

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

JOIN THE NEXTGEN
DATA SCIENCE
ECOSYSTEM
(https://courses.analyticsvidhya.com/courses/computer-learning-version2/?utm_source=home_blog_navbar&utm_medium=stickybanner2)



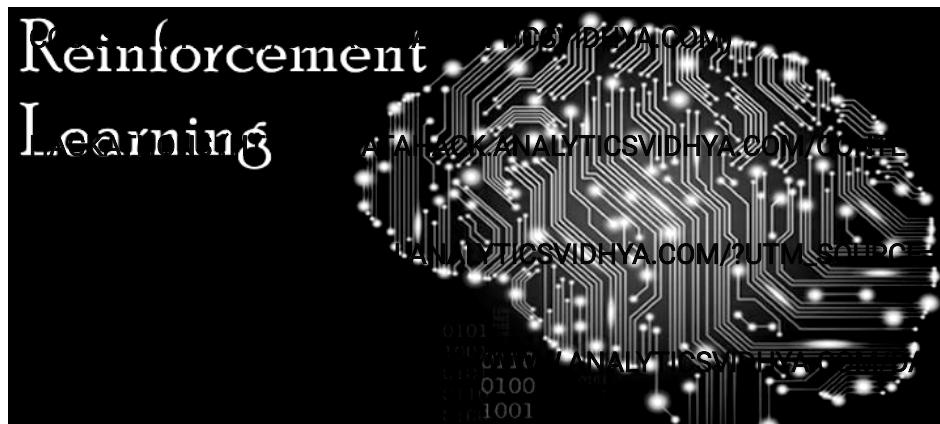
Learn to Solve
Text Classification Problems Using NLP

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=home_blog_navbar&utm_medium=stickybanner2)



Start Building Your First
Computer Vision Model Today

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 ~~CONTACT ([solve problems by *CONTACT \(* interacting with 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.](https://www.analyticsvidhya.com/contact)~~

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.

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

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>)

utm_source=blog&utm_medium=Stickybanner2

Natural Language Processing Using Python

Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

COMPUTER VISION USING DEEP LEARNING

Start Building Your First Computer Vision Model Today

[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR) (https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2?)

Why learn dynamic programming?

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ▾

Apart from being a good starting point for grasping reinforcement

learning, dynamic programming can help find optimal solutions to [HACKATHONS \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL\)](https://datahack.analyticsvidhya.com/contest/all) planning problems faced in the industry, with an important assumption

that the specifics of the environment are known. DP presents a good

~~DATA MINING~~ ([HTTPS://DATAMINING.ANALYTICSVIDHYA.COM/](https://datamining.analyticsvidhya.com/)) ~~?utm_source=HOME_BLOG_NAVBAR~~ complex problems.

[DATAHACK SUMMIT 2019 \(HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?\)](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=HOME_BLOG_NAVBAR) ONLINE COURSE

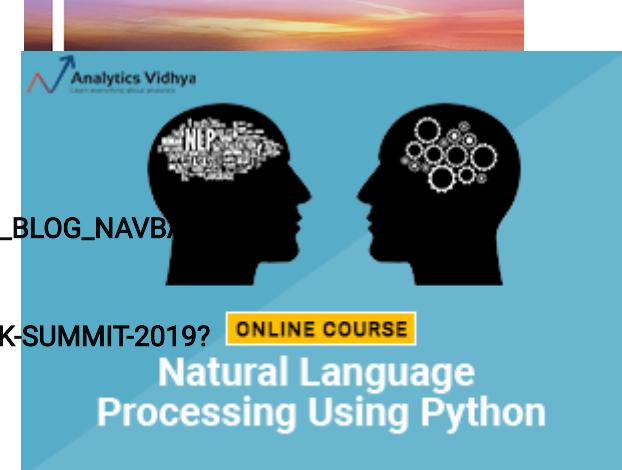
Sunny's Motorbike Rental Company

Sunny manages a motorbike rental company in Ladakh. Being near the [CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT?\)](https://www.analyticsvidhya.com/contact)

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.

[utm_source=blog&utm_medium=Stickybanner2?\)](https://courses.analyticsvidhya.com/courses/natural-language-processing-using-python-with-pytorch?utm_source=blog&utm_medium=Stickybanner2)

[utm_source=blog&utm_medium=Stickybanner2?\)](https://courses.analyticsvidhya.com/courses/natural-language-processing-using-python-with-pytorch?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp?utm_source=blog&utm_medium=Stickybanner2)



Start Building Your First
Computer Vision Model Today



Bikes are rented out for Rs 1200 per day and are available for renting the day after they are returned.

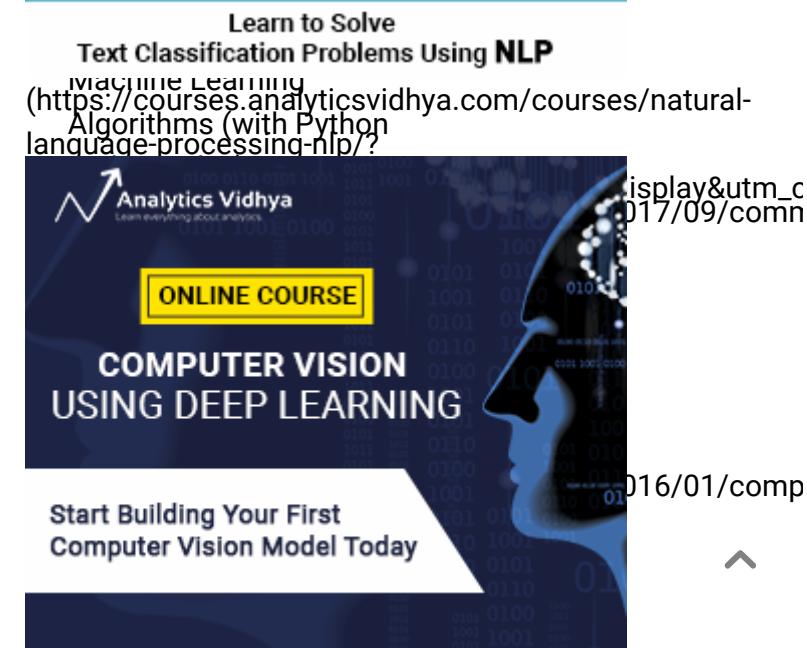
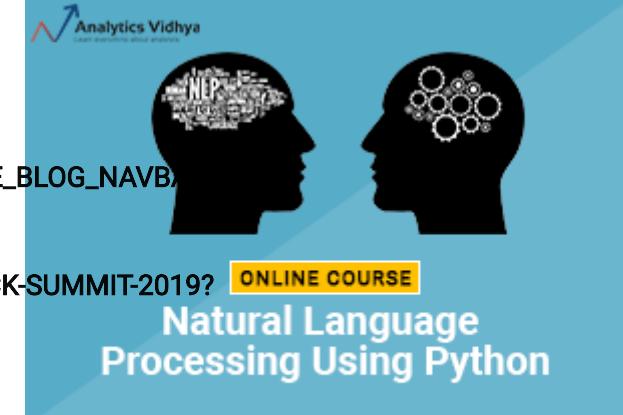
- Sunny can move the bikes from 1 location to another and incurs [CONTACT \(\[HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT?\]\(https://www.analyticsvidhya.com/contact?utm_source=HOME_BLOG_NAVBAR\)\)](https://www.analyticsvidhya.com/contact?utm_source=HOME_BLOG_NAVBAR) a cost of Rs 100
- With experience Sunny 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 $h(n)$

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.

([HTTPS://WWW.ANALYTICSVIDHYA.COM/COURSES/COMPUTER-SCIENCE-DATA-SCIENCE-PROJECTS-TO-BOOST-YOUR-LEARNING-VERSION2/?utm_source=blog&utm_medium=Stickybanner2](https://www.analyticsvidhya.com/courses/computer-science-data-science-project-to-boost-your-learning-version2/?utm_source=blog&utm_medium=Stickybanner2))

can be accessed freely)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2018/05/24-NLP-PROJECTS-TO-BOOST-YOUR-LEARNING-VERSION2/](https://www.analyticsvidhya.com/blog/2018/05/24-nlp-projects-to-boost-your-learning-version2/))



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, then [courses \(https://www.analyticsvidhya.com/courses/reinforcement-learning-dp\)](https://www.analyticsvidhya.com/courses/reinforcement-learning-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 [HACKATHONS \(https://datahack.analyticsvidhya.com/contest/all\)](https://datahack.analyticsvidhya.com/contest/all) key concepts in RL such as policy, reward, value function and more.

[DATAMIN \(https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar\)](https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar)

Table of Contents

[DATAHACK SUMMIT 2019 \(https://www.analyticsvidhya.com/datahack-summit-2019\)](https://www.analyticsvidhya.com/datahack-summit-2019)

1. Understanding Agent-Environment interface using tic-tac-toe
[utm_source=home_blog_navbar Process](https://www.analyticsvidhya.com/courses/reinforcement-learning-dp)

1. Value Function:

1. How good it is to be in a given state?
[CONTACT \(https://www.analyticsvidhya.com/contact/\)](https://www.analyticsvidhya.com/contact)
 2. How good an action is at a particular state?

2. Solving an MDP: The Math

1. Bellman Expectation Equation
2. Bellman Optimality Equation

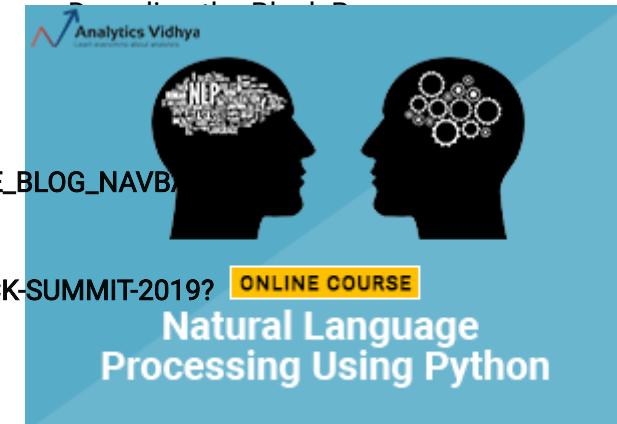
3. Dynamic Programming

1. Policy Iteration

1. Policy Evaluation
2. Policy Improvement

2. Value Iteration

(<https://courses.analyticsvidhya.com/courses/computer-regression-techniques>)
 you should know
[\(https://www.analyticsvidhya.com/blog/2015/08/computer-regression/\)](https://www.analyticsvidhya.com/blog/medium=stickybanner?utm_source=blog&utm_medium=Stickybanner2/)



Learn to Solve
 Text Classification Problems Using **NLP**

Innovative Data Visualizations you language-processing-nlp/?)



4. DP in action: Finding optimal policy for Frozen Lake environment
BLOG ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR))
 using Python

1. Frozen Lake Environment

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ▾

3. Value Iteration in python

HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

~~Understanding Agent Environment Interface using Tic-tac-toe~~

DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=HOME_BLOG_NAVBAR))

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/TICTAC-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:

[CONTACT](https://www.analyticsvidhya.com/contact) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT](https://www.analyticsvidhya.com/contact))

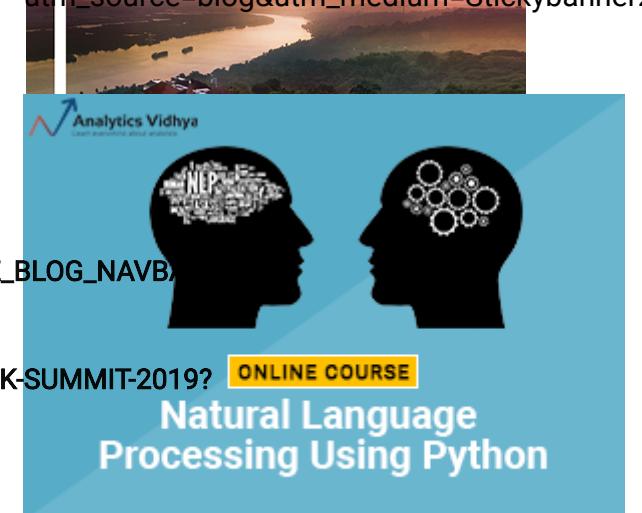
Can you define a rule-based framework to design an efficient bot?

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:

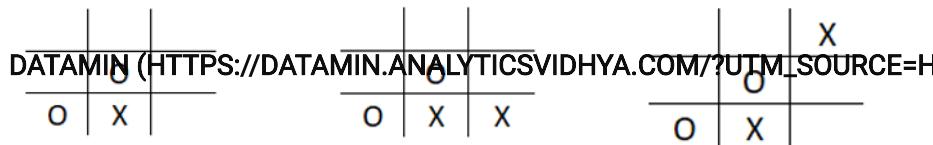
([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/COMPUTER-VISION-USING-DEEP-LEARNING-VERSION2/](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/)?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
 Text Classification Problems Using **NLP**
[\(HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/NATURAL-LANGUAGE PROCESSING-NLP/?\)](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)



- 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 **state** 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**.
- [HACKATHONS](https://datahack.analyticsvidhya.com/contest/all) ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))



[DATAMIN](https://datamin.analyticsvidhya.com/) ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar))

- Once the state is known, the bot must take an **action** to be optimum to win the game (**policy**)
- This move will result in a new scenario with new **CONTACT** ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/)) 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:

RECENT POSTS analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/

?utm_source=blog&utm_medium=Stickybanner2)

Everything you Should Know about p-value from

Analytics Vidhya
Learn everything about analytics

NLP

ONLINE COURSE

Natural Language Processing Using Python

Learn to Solve Text Classification Problems Using **NLP** (<https://www.analyticsvidhya.com/courses/natural-language-processing/2019/09/feature-engineering-and-nlp>)

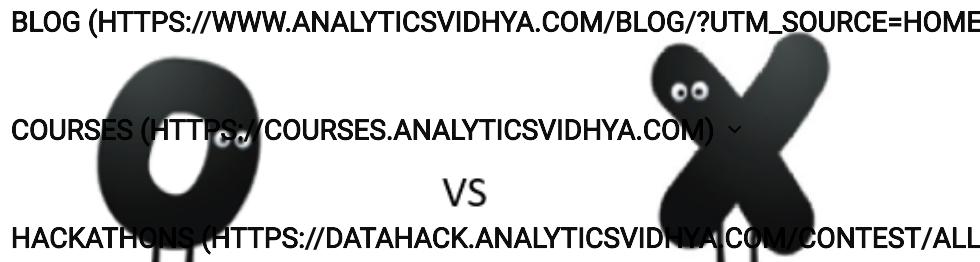
?display&utm_c

Analytics Vidhya
Learn everything about analytics

ONLINE COURSE

COMPUTER VISION USING DEEP LEARNING

Start Building Your First Computer Vision Model Today



DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR](https://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR))

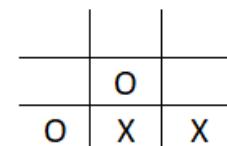
Consider the following situation encountered in tic-tac-toe:

DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR](https://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR))



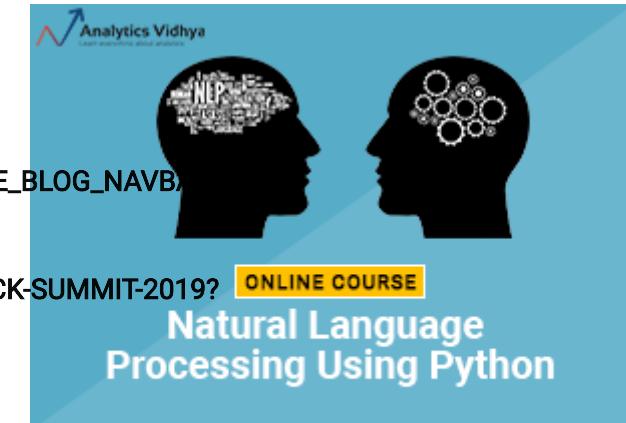
CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://WWW.ANALYTICSVIDHYA.COM/CONTACT/))

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



https://courses.analyticsvidhya.com/courses/computer-step-deep-learning-tutorial-video-classification-python?utm_medium=Stickybanner2

SEPTEMBER 3, 2019



Learn to Solve
Text Classification Problems Using **NLP**

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_medium=Stickybanner2)



ONLINE COURSE

**COMPUTER VISION
USING DEEP LEARNING**

Start Building Your First
Computer Vision Model Today

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 [penalty](https://www.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/) to X for this action in the next trial. We say that this action in the given state would correspond to a negative reward and should not be considered as an optimal action in this situation.

[DATAHACK SUMMIT 2019 \(HTTPS://WWW.ANALYTICSVIDHYA.COM/UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/courses/natural-language-processing-using-python/) winning in the next move:



[CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/)

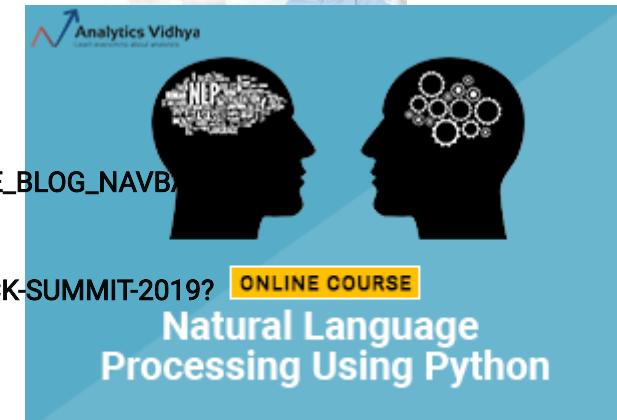
Introduction to Markov Decision Process

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:

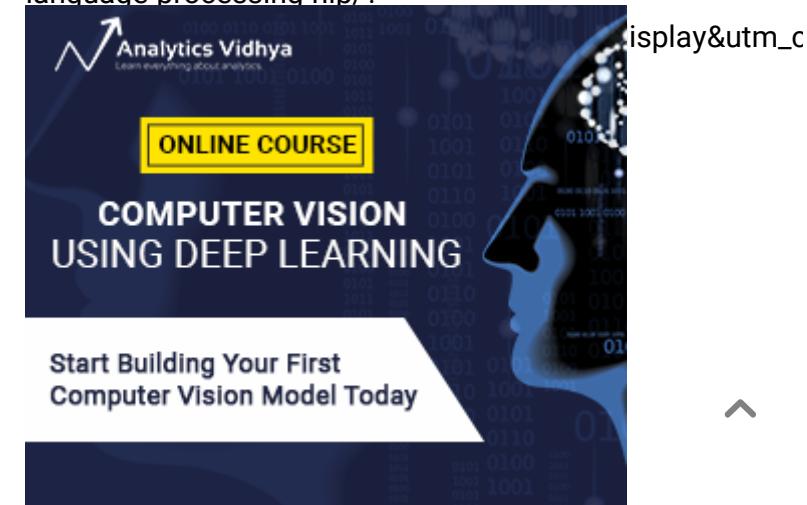
- A set of possible world states S
- A set of possible actions A
- A real valued reward function R(s,a)
- A description T of each action's effects in each state

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?>)



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

[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR)

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>)

Any random process in which the probability of being in a given state

depends only on the previous state is called a markov process.

[COURSES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/COURSES\)](https://www.analyticsvidhya.com/courses)

In other words, in the markov decision process setup, the

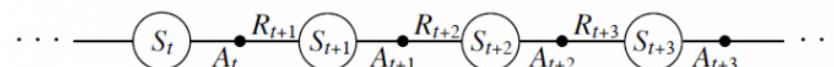
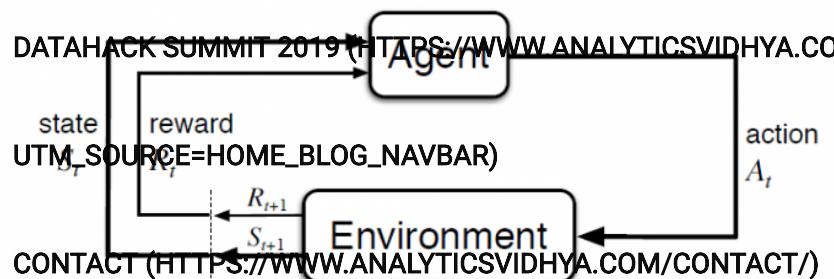
[HACKATHONS \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL\)](https://datahack.analyticsvidhya.com/contest/all)

environment's response at time $t+1$ depends only on the state and

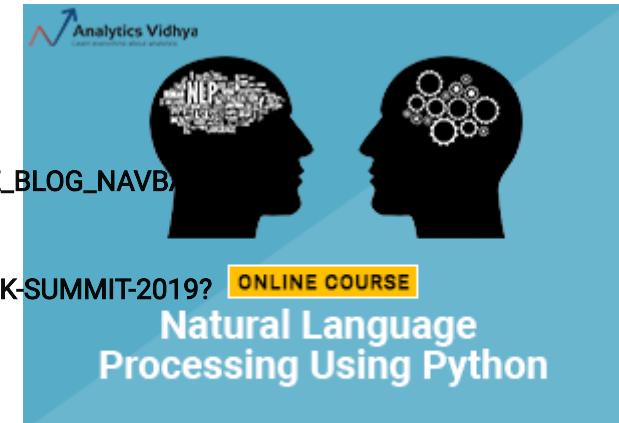
action representations at time t , and is independent of whatever

happened in the past.

[DATAMIN \(HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://datamin.analyticsvidhya.com/?utm_source=HOME_BLOG_NAVBAR)



- St: State of the agent at time t
- At: Action taken by agent at time t
- Rt: Reward obtained at time t



(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



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 course (<https://www.analyticsvidhya.com>) gives in the long run. Total reward at any time instant t is given by:

HACKATHONS (<https://datahack.analyticsvidhya.com/contest/all>)

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

DATAMIN (https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar) 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 γ as a discounting factor and each reward is discounted by this factor as follows:

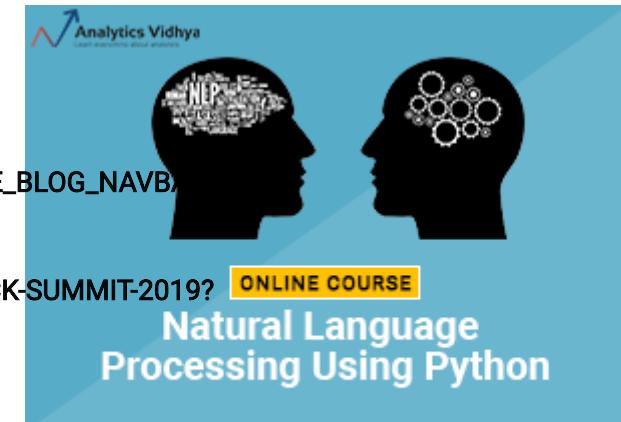
CONTACT (<https://www.analyticsvidhya.com/contact/>)

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

For discount factor < 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 (γ close to 1) or short term (γ close to 0).

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

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?>)



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 π represents how good a state is for a agent to be in it. How do we calculate the value function? What is the average reward that the agent will get starting from the current state under policy π ?

HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

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

DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar))

E in the above equation represents the expected reward at each state

if the agent follows policy π and S represents the set of all possible states.

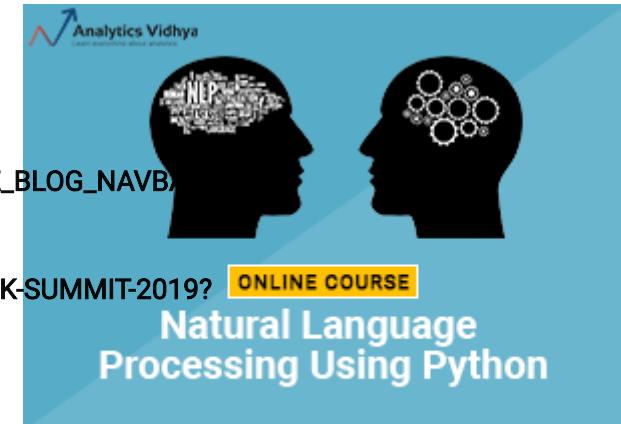
Policy, as discussed earlier, is the mapping of probabilities of taking each possible action at each state ($\pi(a|s)$). The policy might also be denoted in terms of values exactly like a value function, but does not give probabilities.

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?

([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/COMPUTER-VISION-USING-DEEP-LEARNING-VERSION2/](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/)?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/NATURAL-LANGUAGE PROCESSING-NLP/](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)?)



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 function as:

$$\text{HACKATHONS} (\text{HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL}) \quad \left[\sum_{k=0}^{\infty} \right]$$

DATAMIN (https://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR)
This is the expected return the agent will get if it takes action At at time t, given state St, and thereafter follows policy π .

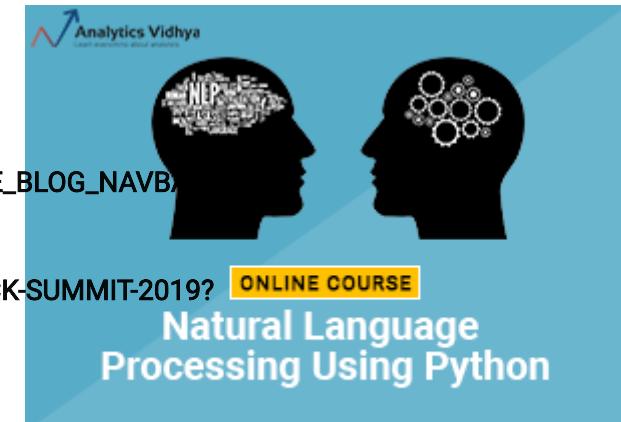
DATAHACK SUMMIT 2019 (https://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR)

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

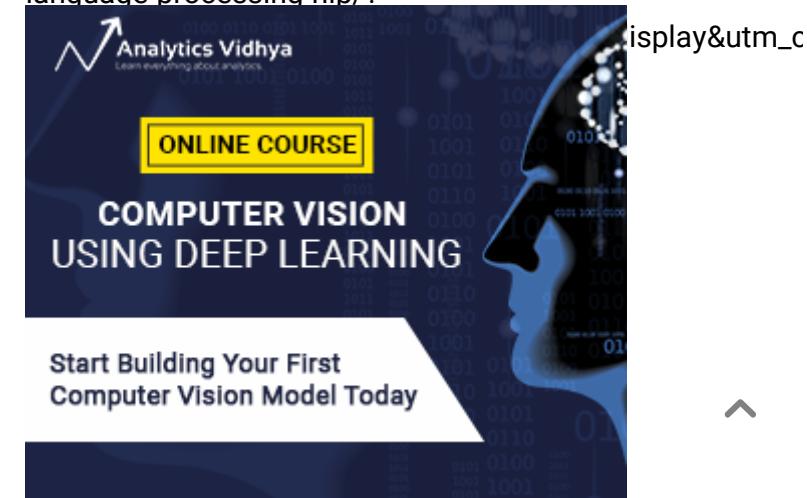
Bellman was an applied mathematician who derived equations that help to solve an Markov Decision Process.

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

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**
(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blog&utm_medium=Stickybanner2)



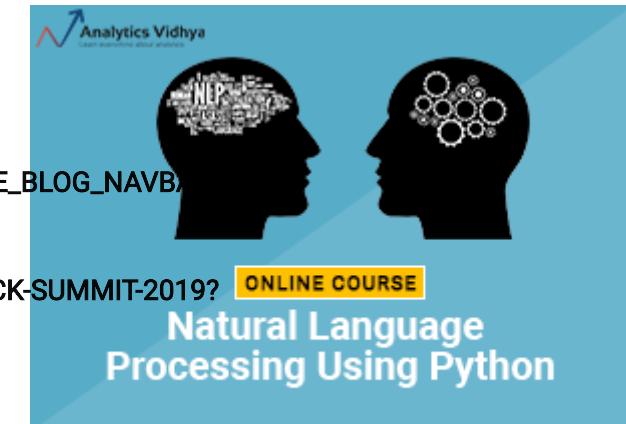
$$\begin{aligned}
 v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\
 &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
 &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi}[G_{t+1} \mid S_{t+1} = s']] \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')], \quad \text{for all } s \in \mathcal{S}, \\
 &\text{HACKATHONS } (\text{HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL})
 \end{aligned}$$

In this equation we have the value function for a given policy π represented in terms of the value function of the next state.

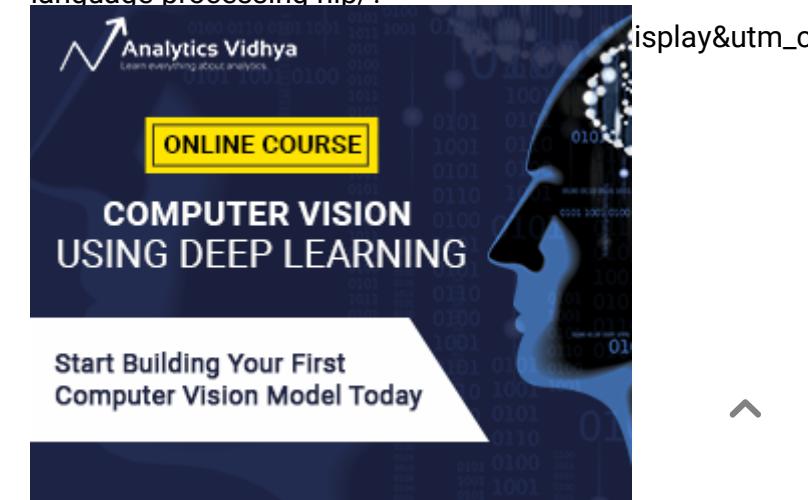
Choose an action a , with probability $\pi(a|s)$ at the state s , which leads to state s' with prob $p(s'|s,a)$. This gives a reward $[r + \gamma * v_{\pi}(s)]$ as given in the equation above.

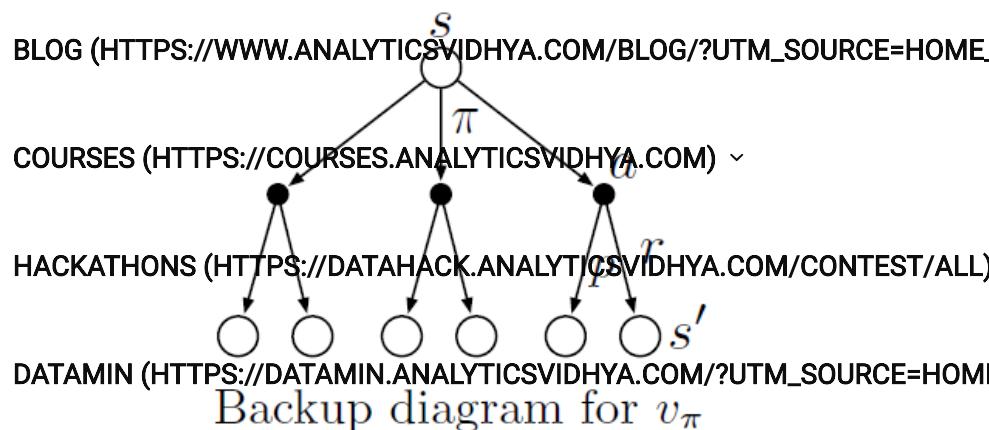
This is called the Bellman Expectation Equation. The value information from successor states is being transferred back to the current state, and this can be represented efficiently by something called a backup diagram as shown below.

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=HOME_BLOG_NAVBAR)





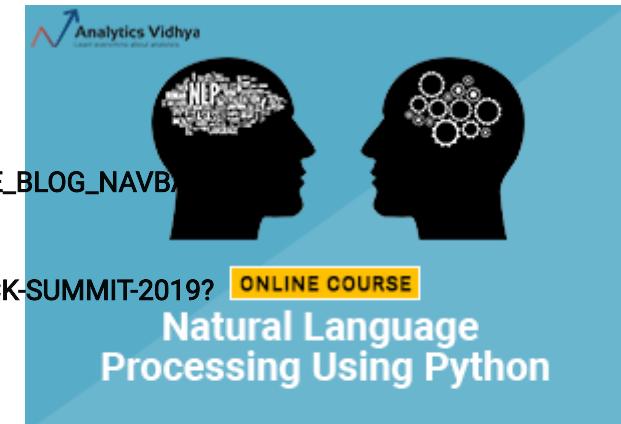
DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=HOME_BLOG_NAVBAR))
The Bellman expectation equation averages over all the possibilities, weighting each by its probability of occurring. 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.

CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))
 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 π , such that for no other π can the agent get a better expected return. We want to find a policy which achieves maximum value for each state.

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve Text Classification Problems Using NLP
 (https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_medium=display&utm_c)



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

Note that we might not get a unique policy, as under any situation COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://COURSES.ANALYTICSVIDHYA.COM)) there can be 2 or more paths that have the same return and are still optimal.

HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL))

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

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

DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR](https://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR))

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

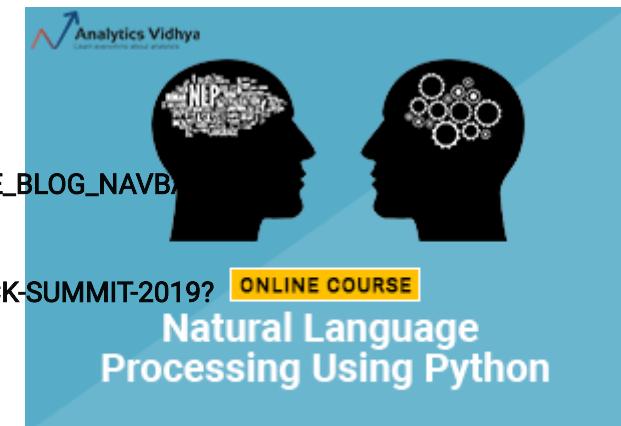
UTM_SOURCE=HOME_BLOG_NAVBAR

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

CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://WWW.ANALYTICSVIDHYA.COM/CONTACT/))

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 π^* , the optimal value function is given by:

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



$$v_{\pi}(s) = \max_{a \in \mathcal{A}(s)} q_{\pi}(s, a)$$

$$= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a]$$

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ~

$$= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

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

DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM](https://datamin.analyticsvidhya.com)) ?UTM_SOURCE=HOME_BLOG_NAVBAR

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

DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019](https://www.analyticsvidhya.com/datahack-summit-2019))?ONLINE COURSE

$$\pi'(s) \doteq \arg \max q_{\pi}(s, a)$$

?UTM_SOURCE=HOME_BLOG_NAVBAR

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

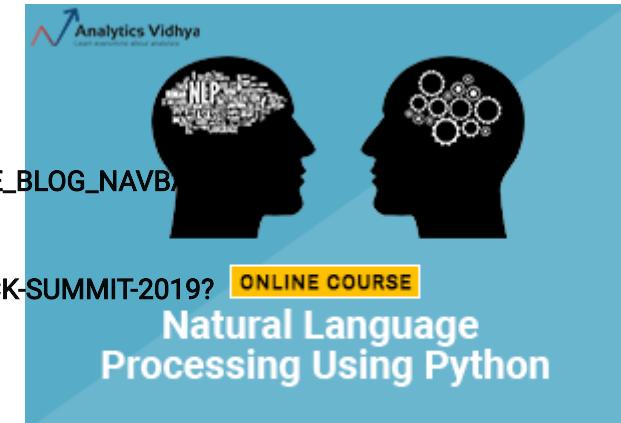
CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT](https://www.analyticsvidhya.com/contact))

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. Given the complete model and specifications of the environment (MDP), we can successfully find an

([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/COMPUTER-VISION-USING-DEEP-LEARNING-VERSION2/](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/)?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/NATURAL-LANGUAGE PROCESSING-NLP/](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/))



Start Building Your First
Computer Vision Model Today

optimal policy for the agent to follow. It contains two main steps:

[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR)

1. Break the problem into subproblems and solve it

[COURSES \(HTTPS://COURSES.ANALYTICSVIDHYA.COM\)](https://courses.analyticsvidhya.com/)

Solutions to subproblems recorded or stored for reuse to find overall optimal solution to the problem at hand

[HACKATHONS \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL\)](https://datahack.analyticsvidhya.com/contest/all)

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

[DATAMIN \(HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://datamin.analyticsvidhya.com/?utm_source=HOME_BLOG_NAVBAR)

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

[DATAHACK SUMMIT 2019 \(HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019\)](https://www.analyticsvidhya.com/datahack-summit-2019)

[UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/contact/?utm_source=HOME_BLOG_NAVBAR)

Policy Evaluation: Find out how good a policy is?

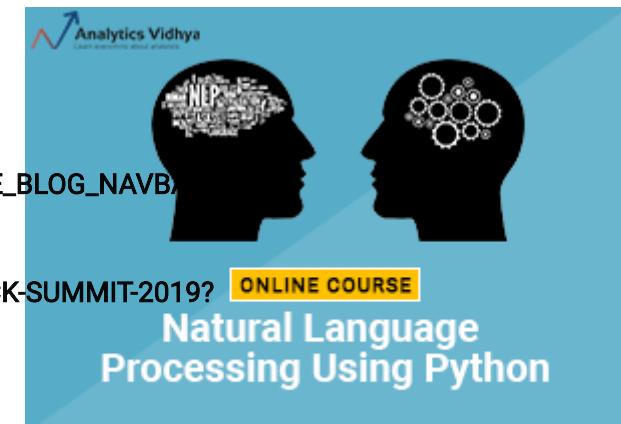
[CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/)

Policy evaluation answers the question of how good a policy is. Given an MDP and an arbitrary policy π , 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.

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

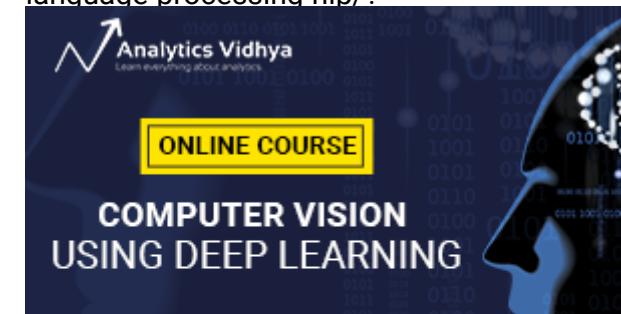
To produce each successive approximation v_{k+1} from v_k , iterative policy evaluation applies the same operation to each state s . It replaces the old value of s with a new value obtained from the old values of the successor states of s , and the expected immediate

[https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2\)](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

[https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blog&utm_medium=Stickybanner2\)](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blog&utm_medium=Stickybanner2)

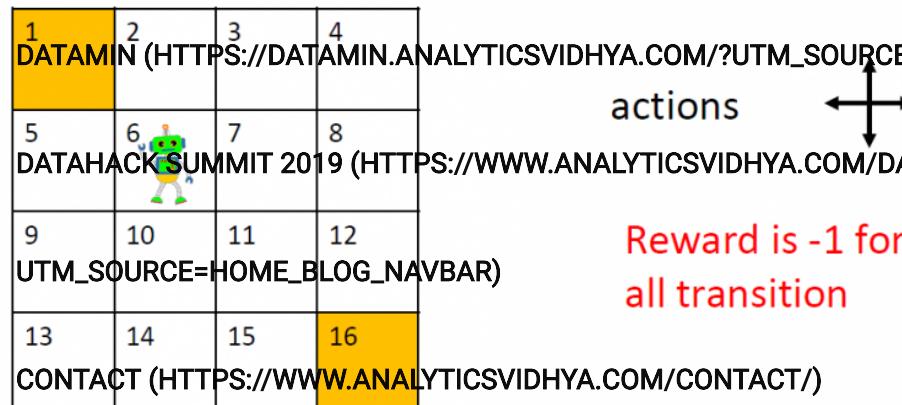


rewards, along all the one-step transitions possible under the policy being evaluated, until it converges to the true value function of a given policy π .

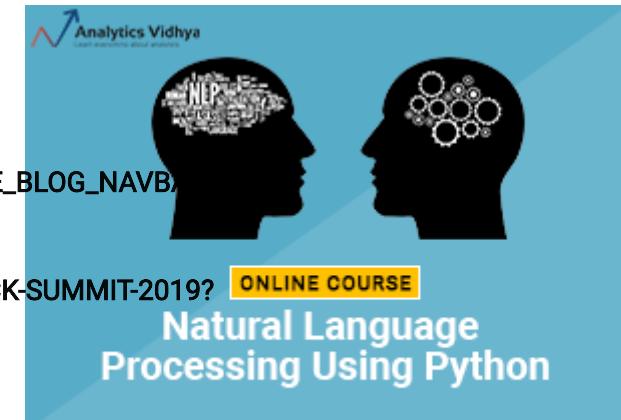
(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ~
Let us understand policy evaluation using the very popular example of Gridworld.

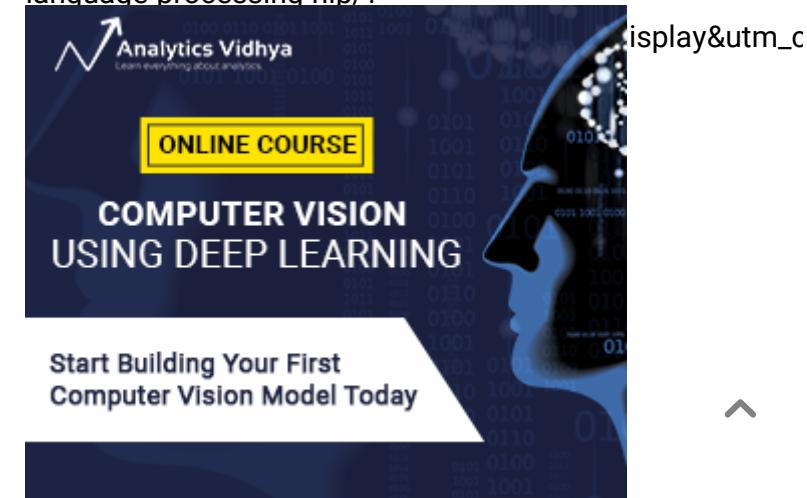
HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)



A bot is required to traverse a grid of 4×4 dimensions to reach its goal (1 or 16). 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,...,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 initialising v_0 for the random policy to all 0s.



Learn to Solve Text Classification Problems Using **NLP**
(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR)	0.0	0.0	0.0	0.0
COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM)	0.0	0.0	0.0	0.0
HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)	0.0	0.0	0.0	0.0
DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR)	0.0	0.0	0.0	0.0

DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?)

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

6:

UTM_SOURCE=HOME_BLOG_NAVBAR)

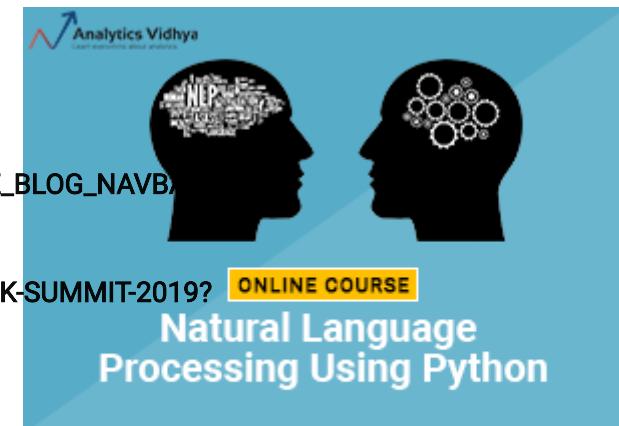
1	2	3	4
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)			
5	6	7	8
9	10	11	12
13	14	15	16

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

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

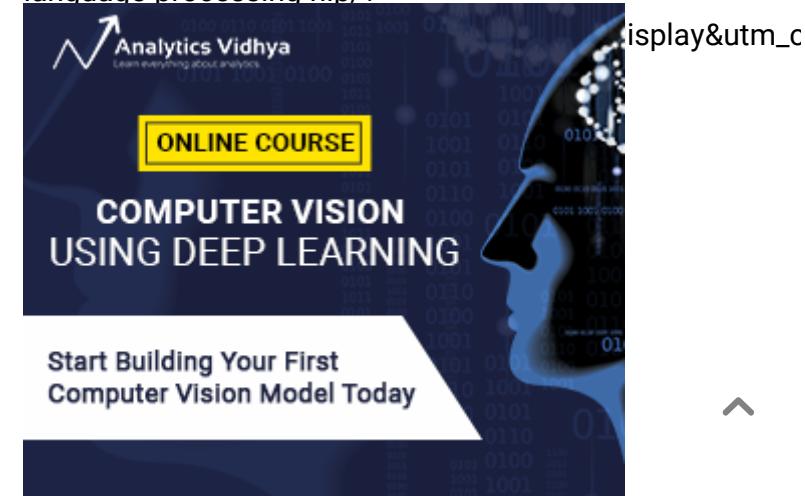
For terminal states $p(s'|s, a) = 0$ and hence $v_k(1) = v_k(16) = 0$ for all k.
So v_1 for the random policy is given by:

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?)



ONLINE COURSE
COMPUTER VISION
USING DEEP LEARNING

Start Building Your First
Computer Vision Model Today

BLOG (HTTP://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR)	0.0	-1.0	-1.0	-1.0
COURSES (HTTP://COURSES.ANALYTICSVIDHYA.COM)	-1.0	1.0	-1.0	-1.0
HACKATHONS (HTTP://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)	-1.0	-1.0	1.0	-1.0
DATAMIN (HTTP://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR)	-1.0	-1.0	-1.0	1.0

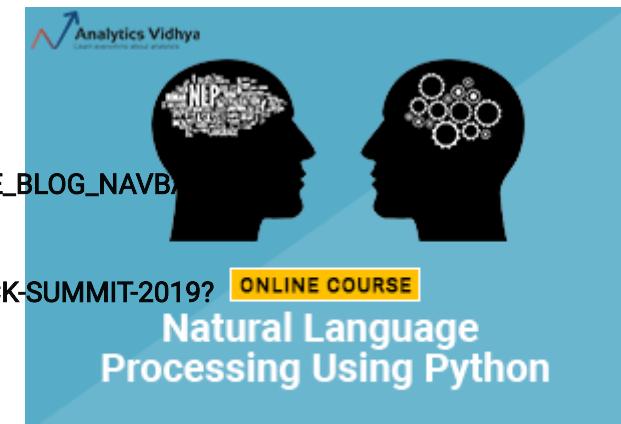
DATAHACK SUMMIT 2019 (HTTP://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?)

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

UTM_SOURCE=HOME_BLOG_NAVBAR)

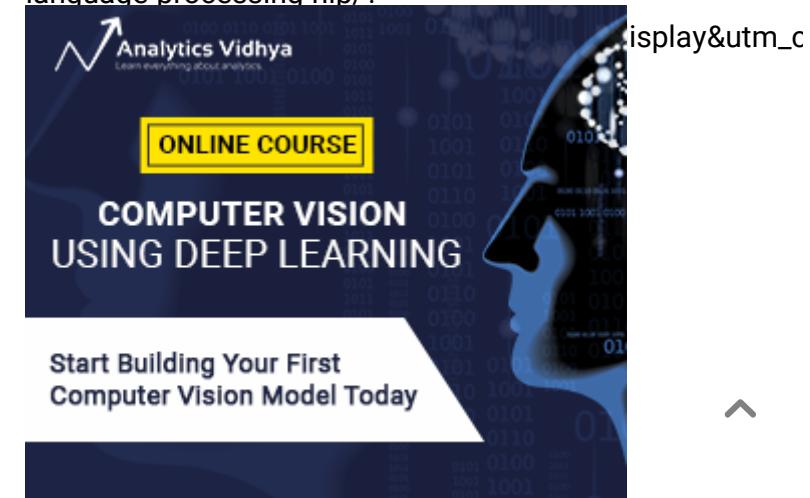
CONTACT (HTTP://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?)



$$\text{BLOG}(\text{HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm_source=HOME_BLOG_NAVBAR}) = 0.25 \sum_{a \in \{u, d, l, r\}} \sum_{s'} p(s'|s, a) \text{value}_{t+1}(s')$$

$$= -1 = \begin{cases} -1, s' \in S \\ 0, s' \in S^+ \setminus S \end{cases}$$

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com))

$$= 0.25 * \{p(2|6, u)[-1 - \gamma] + p(10|6, d)[-1 - \gamma]$$

$$+ p(5|6, l)[-1 - \gamma] + p(7|6, r)[-1 - \gamma]\}$$

$$\gamma = 1$$

$$= 0.25 * \{-2 - 2 - 2 - 2\}$$

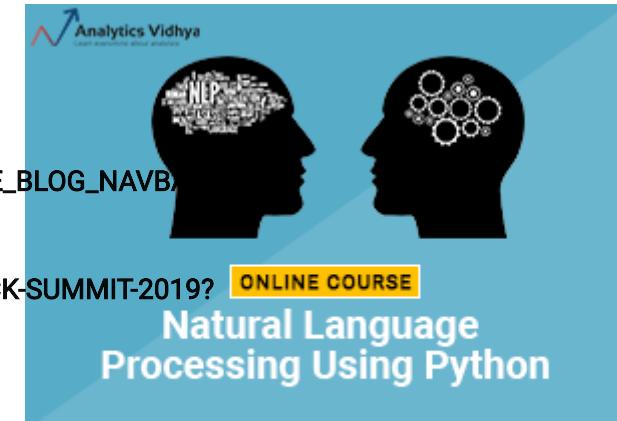
$$\text{DATAMIN}(\text{HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR})$$

DATAHACK SUMMIT	2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019)			
UTM_SOURCE=HOME_BLOG_NAVBAR			8	
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT)	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:

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>?)



Start Building Your First
Computer Vision Model Today

$$\text{BLOG}(\text{HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR}) = 0.25 \sum_{a \in \{u,d,l,r\}} \sum_{s'} p(a|s) v_2(s') = -1$$

$\begin{cases} -1, s' \in S \\ 0, s' \in S^+ \setminus S \end{cases}$

COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM)

$$= 0.25 * \{p(2|2,u)[-1 - \gamma] + p(6|2,d)[-1 - \gamma]$$

$$+ p(1|2,l)[-1 - \gamma * 0] + p(3|2,r)[-1 - \gamma]\}$$

$$\stackrel{\gamma = 1}{=} 0.25 * \{-2 - 2 - 1 - 2\}$$

DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR)

$$\Rightarrow v_2(2) = -1.75$$

DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=HOME_BLOG_NAVBAR)

UTM_SOURCE=HOME_BLOG_NAVBAR

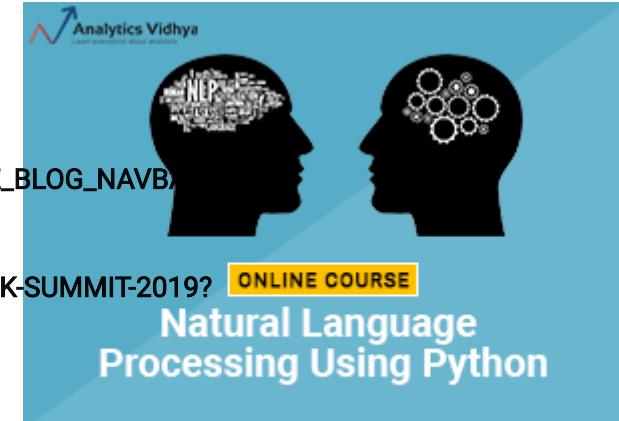
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

$\Rightarrow v_2$ for the random policy:

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

If we repeat this step several times, we get v_{π_1} :

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

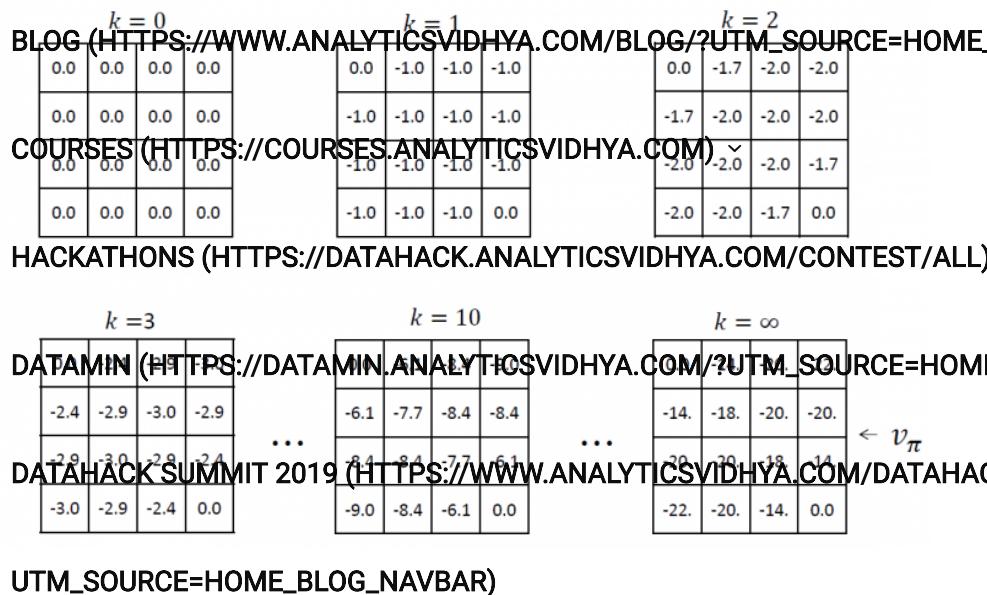


Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



Start Building Your First
Computer Vision Model Today

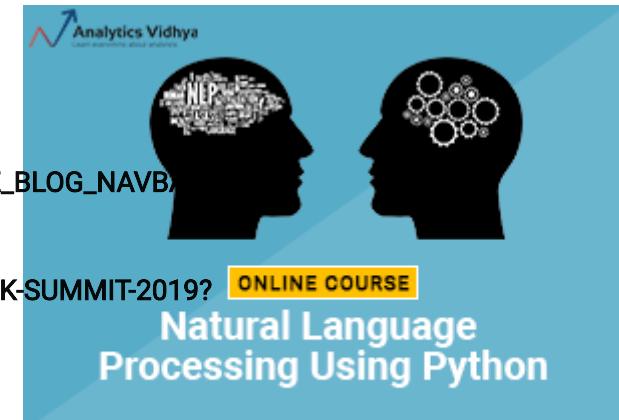


Policy Improvement: Improve an arbitrary policy

Using policy evaluation we have determined the value function v for an arbitrary policy π . 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 π . Let's say we select a in s , and after that we follow the original policy π . The value of this way of behaving is represented as:

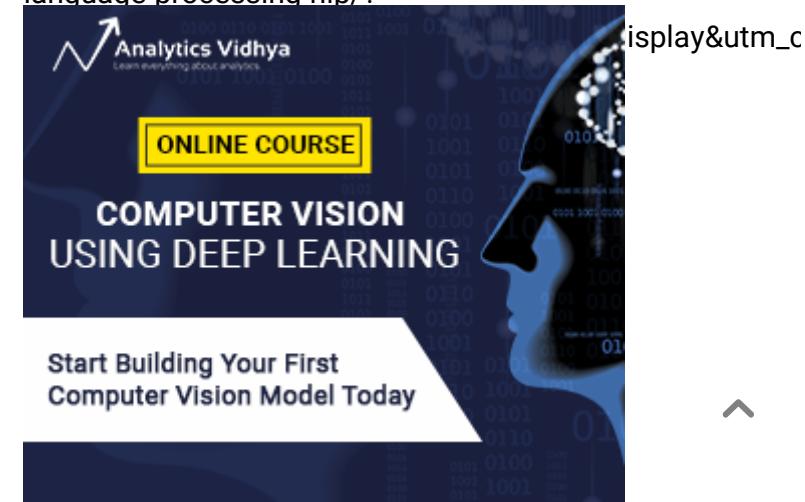
$$\begin{aligned} q_\pi(s, a) &\doteq \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]. \end{aligned}$$

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

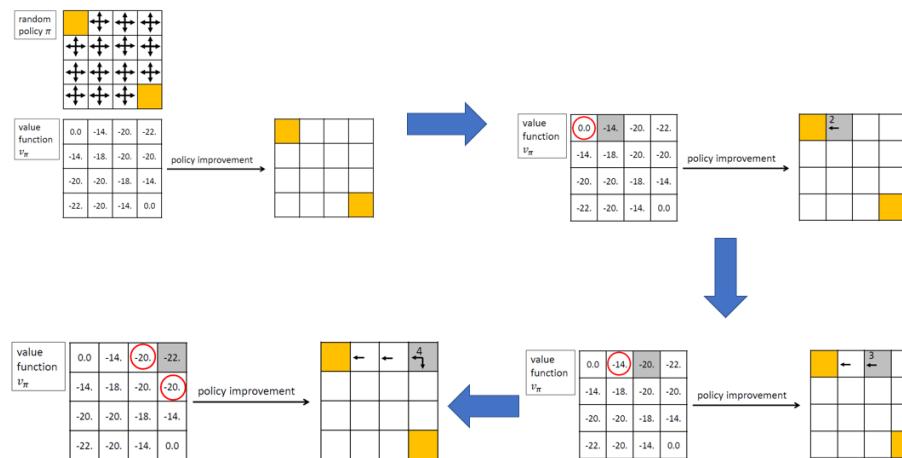
(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



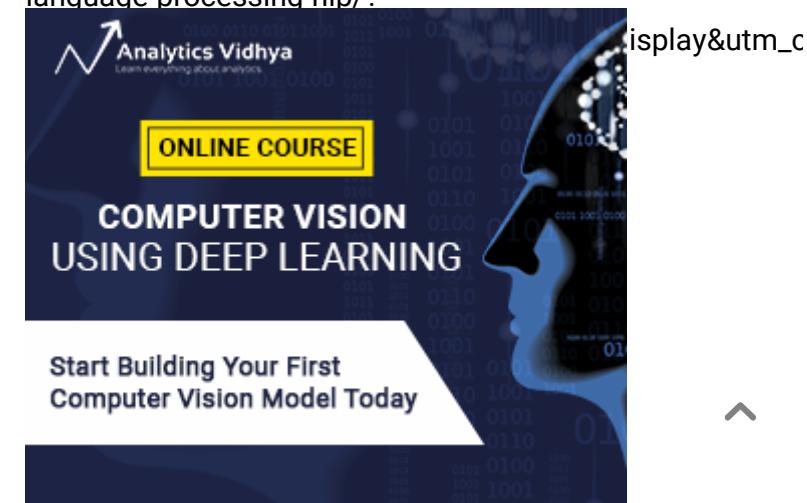
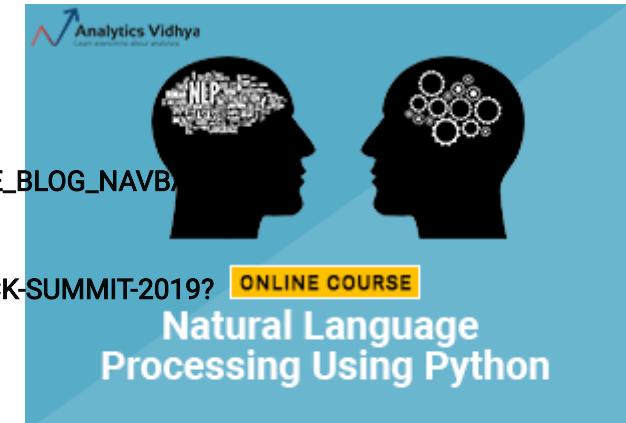
If this happens to be greater than the value function $v_{\pi}(s)$, it implies that the new policy π' would be better to take. We do this iteratively for all states to find the best policy. Note that in this case, the agent would be looking one step ahead. The steps involved are:

HACKATHONS (<https://datahack.analyticsvidhya.com/contest/all>)
 Let's get back to our example of gridworld. Using v_{π} , the value function obtained for random policy π , we can improve upon π by following the **POLICY IMPROVEMENT** (<https://datahack.analyticsvidhya.com/policy-improvement/>). For 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.

AUTOMATION (<https://www.analyticsvidhya.com/home/blog/navbar/>)
 The 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 **CONTACT** (<https://www.analyticsvidhya.com/contact/>) policy.



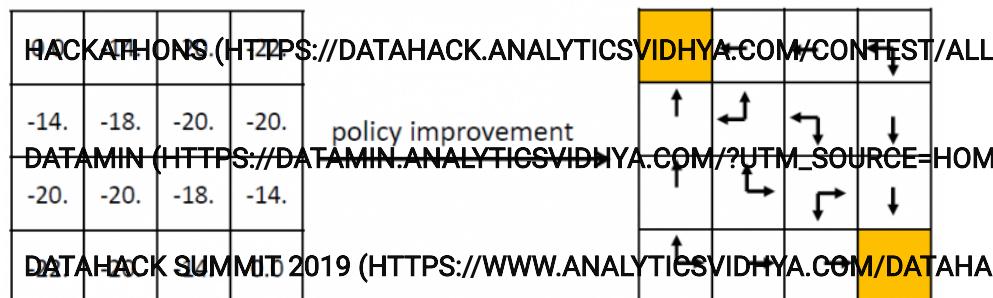
(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>)
 utm_source=blog&utm_medium=Stickybanner2)



[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR) (https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

Overall, after the policy improvement step using v_{π} , we get the new policy π'

[COURSES \(HTTPS://COURSES.ANALYTICSVIDHYA.COM\)](https://courses.analyticsvidhya.com) ▾

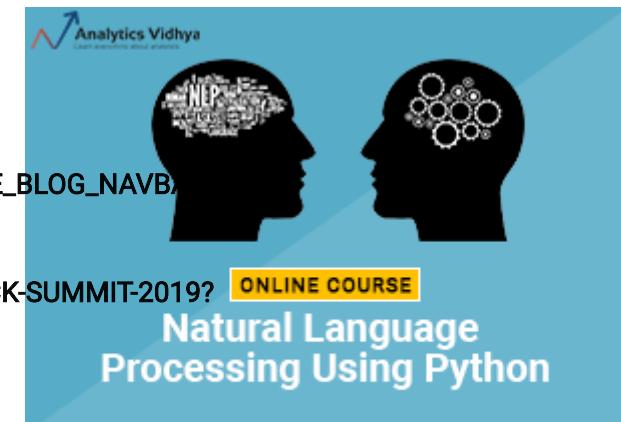


[UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/blog/navbar)

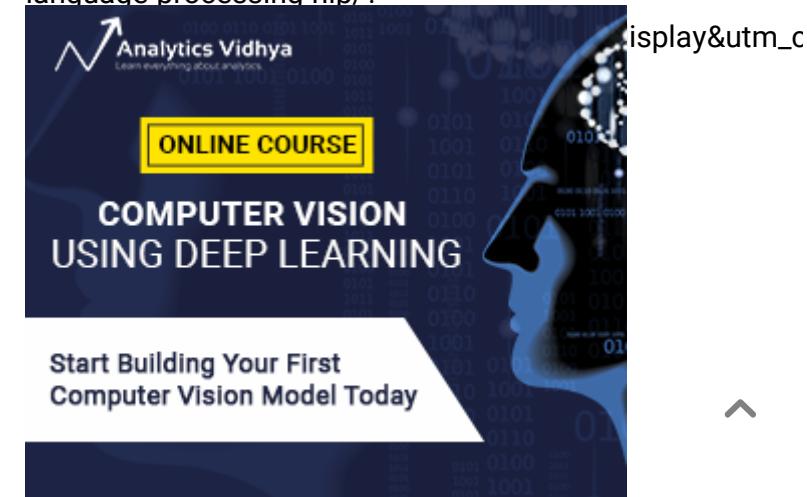
Looking at the new policy, it is clear that it's much better than the random policy. However, we should calculate v_{π}' using the policy evaluation technique we discussed earlier to verify this point and for [CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact) better understanding.

Policy Iteration: Policy Evaluation + Policy Improvement

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



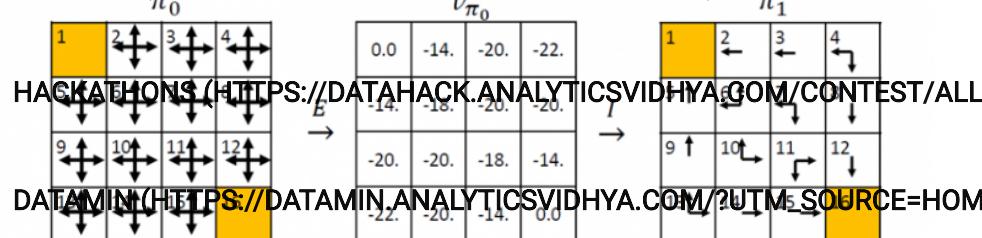
[Learn to Solve Text Classification Problems Using NLP](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)
(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blog&utm_medium=Stickybanner2)



BLOG ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR))

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com))



HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

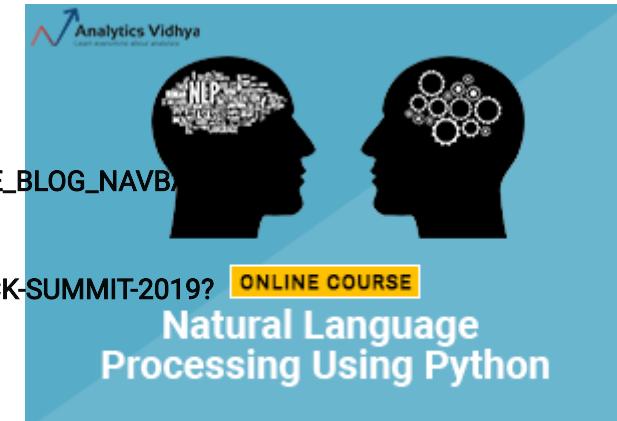
DATA MINING ([HTTPS://DATAMINING.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR](https://datamining.analyticsvidhya.com/?utm_source=HOME_BLOG_NAVBAR))

DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=HOME_BLOG_NAVBAR))

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 source of biasing but here 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_{π_t} , we could stop earlier.



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



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 (<https://www.analyticsvidhya.com/courses/analyticsvidhya.com>) ~

<https://www.analyticsvidhya.com/courses/analyticsvidhya.com>

This is similar to https://www.analyticsvidhya.com/courses/analyticsvidhya.com?utm_source=HOME_BLOG_NAVBAR, 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.

UTM_SOURCE=HOME_BLOG_NAVBAR

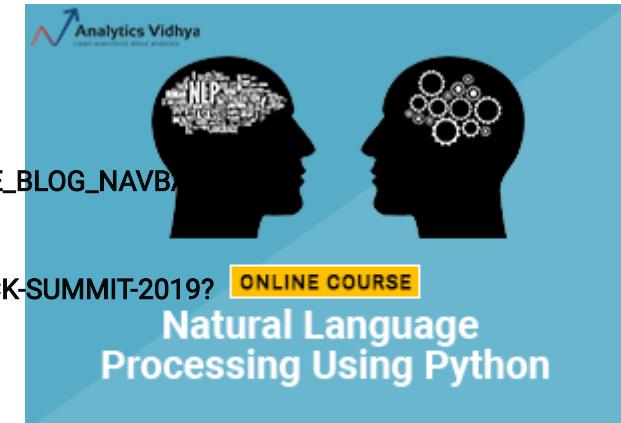
Some important points related to DP:

[CONTACT \(https://www.analyticsvidhya.com/contact/\)](https://www.analyticsvidhya.com/courses/analyticsvidhya.com?utm_source=HOME_BLOG_NAVBAR)

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

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



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 reinforcement learning algorithms. To illustrate dynamic programming here, we will use it to navigate the Frozen Lake environment.

[HACKATHONS \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL\)](https://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)

DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR) Frozen Lake Environment

[DATAHACK SUMMIT 2019 \(HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019\)](https://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019)
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 hole. The 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.

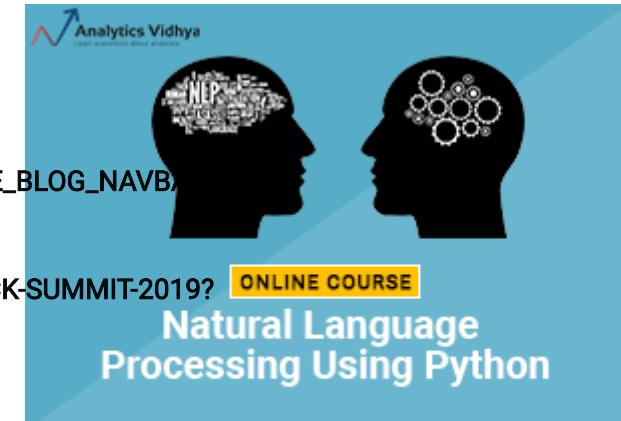
[CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

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

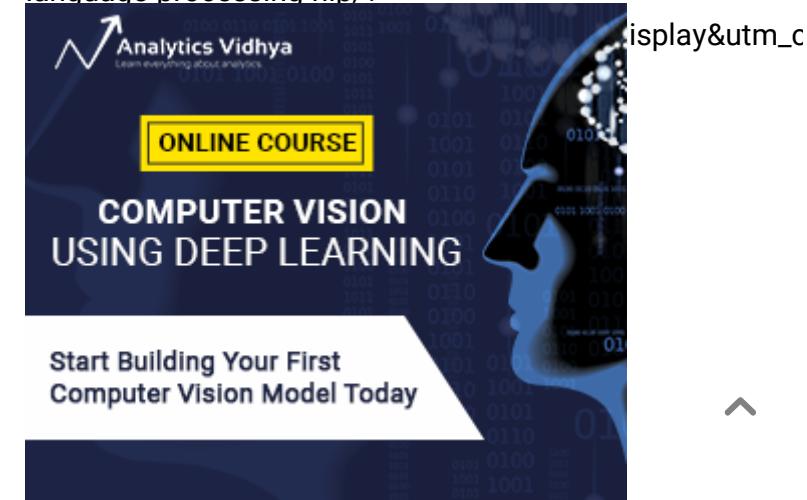
(S: starting point, safe), (F: frozen surface, safe), (H: hole, fall to your doom), (G: goal)

[\(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?\)](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/)
utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

[\(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)



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://openai.com/gym/>)

Once gym library is installed, you can just open a jupyter notebook to [HACKATHONS](#) (<https://datahack.analyticsvidhya.com/contest/all>) get started.

[DATAMIN](#) (https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar)

```
import gym
import numpy as np
```

[DATAHACK SUMMIT 2019](#) (<https://www.analyticsvidhya.com/datahack-summit-2019>)

[UTM_SOURCE=HOME_BLOG_NAVBAR](#)

Now, the env variable contains all the information regarding the frozen lake environment. Before we move on, we need to understand what an [episode](#) (<https://www.analyticsvidhya.com/contact-us>) is used 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

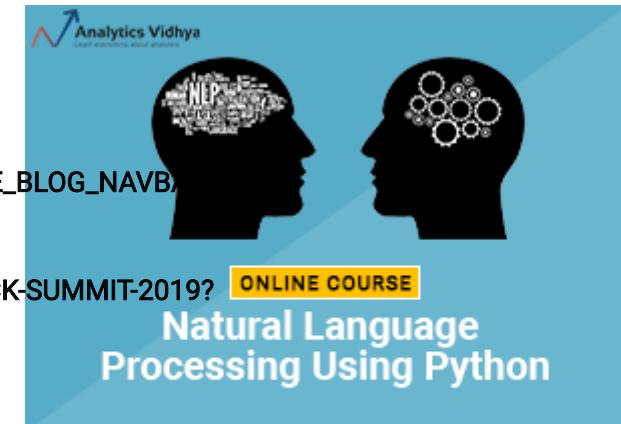
Description of parameters for policy iteration function

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

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>)

utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



discount_factor: MDP discount factor

[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR)

theta: A threshold of a value function change. Once the update to value

function is below this number

[COURSES \(HTTPS://COURSES.ANALYTICSVIDHYA.COM\)](https://courses.analyticsvidhya.com) ▾

max_iterations: Maximum number of iterations to avoid letting the
[HACKATHONS \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL\)](https://datahack.analyticsvidhya.com/contest/all)
 program run indefinitely

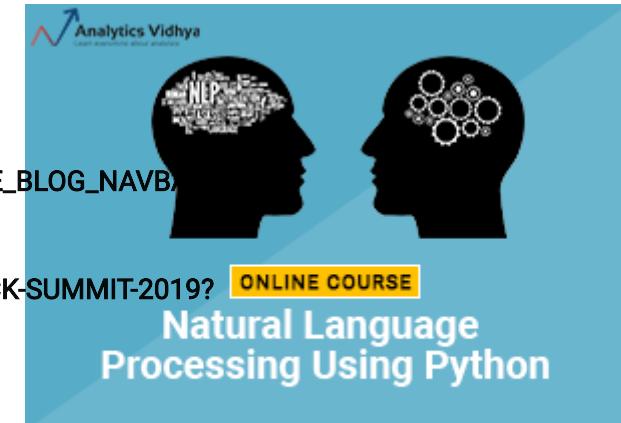
This function will return a vector of sizes, which represent a value
[DATAMINING \(HTTPS://DATAMINING.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://datamining.analyticsvidhya.com/?utm_source=HOME_BLOG_NAVBAR)

[DATAHACK SUMMIT 2019 \(HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=HOME_BLOG_NAVBAR)
 Let's start with the policy evaluation step. The objective is to converge
 to the true value function for a given policy π . We will define a function
[that takes the policy and value function.](https://www.analyticsvidhya.com/policy-evaluation?utm_source=HOME_BLOG_NAVBAR)

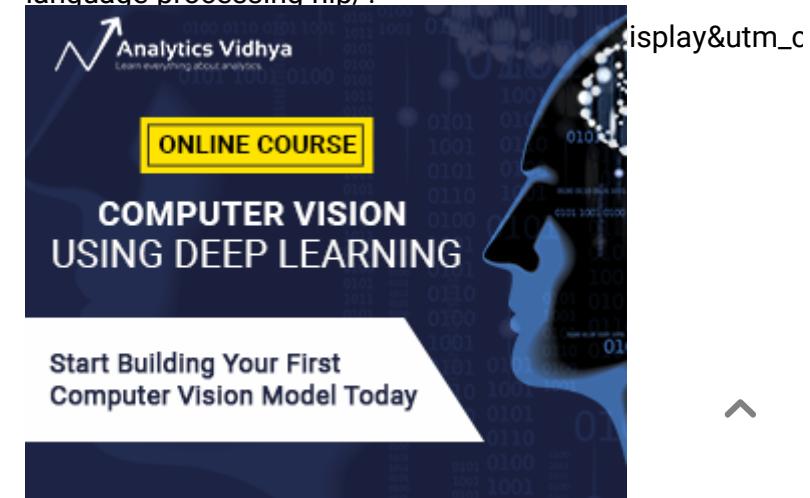
[CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/)

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

utm_source=blog&utm_medium=Stickybanner2)



(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?>)

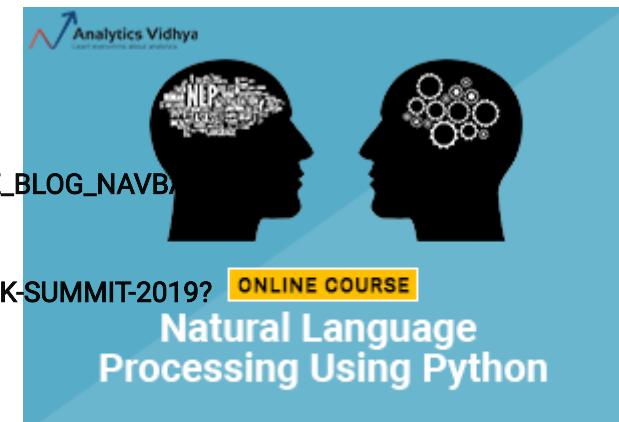


```

BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR)
0, theta=1e-9, max_iterations=1e9):
COURSES# Number of evaluation iterations
evaluation_iterations = 1
HACKATHONS# Initialize a value function for each state as zero
V = np.zeros(environment.nS)
DATAMIN# Repeat until change in value is below the threshold
for i in range(int(max_iterations)):
    # Initialize a change of value function as zero
    DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=home_blog_navbar)
    delta = 0
    UTM_SOURCE=HOME_BLOG_NAVBAR
    # Iterate though each state
    for state in range(environment.nS):
        CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)
        # 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_stat
            e, reward, terminated in environment.P[state][action]:
                # Calculate the expected val

```

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>?utm_source=blog&utm_medium=Stickybanner2)



(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>?display&utm_c



```

ue
BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR)
    v += action_probability * st
    ate_probability * (reward + discount_factor * V[next_state])
COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ~

        # Calculate the absolute change of value
HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)
e function

    delta = max(delta, np.abs(V[state] - V))

        # Update value function
DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?V[state]
evaluation_iterations += 1
UTM_SOURCE=HOME_BLOG_NAVBAR)

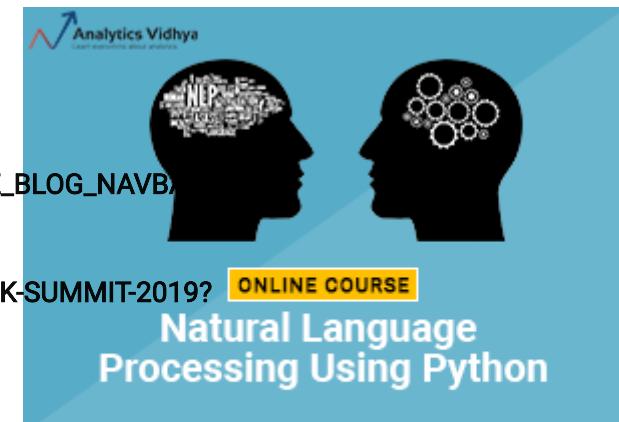
        # Terminate if value change is insignificant
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

    print(f'Policy evaluated in {evaluatio
n_iterations} iterations.')
return V

```

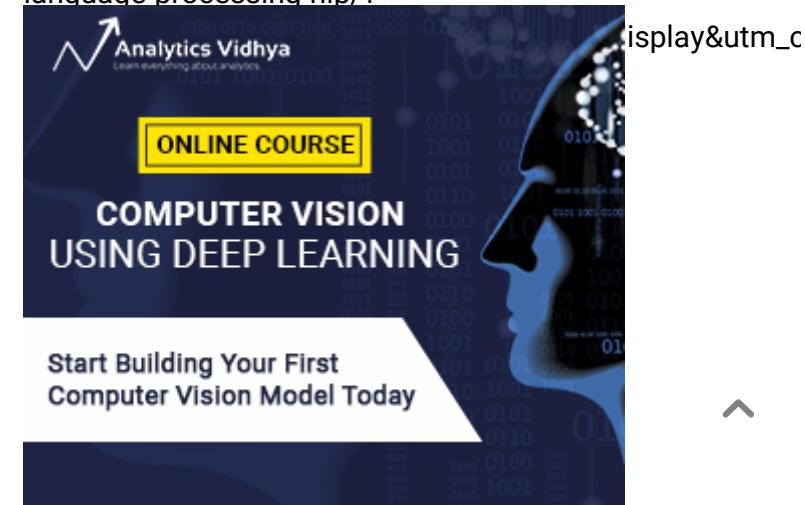
Now coming to the **policy improvement** part of the policy iteration algorithm. We need a helper function that does one step lookahead to calculate the state-value function. This will return an array of length n_A containing expected value of each action

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?>)



```

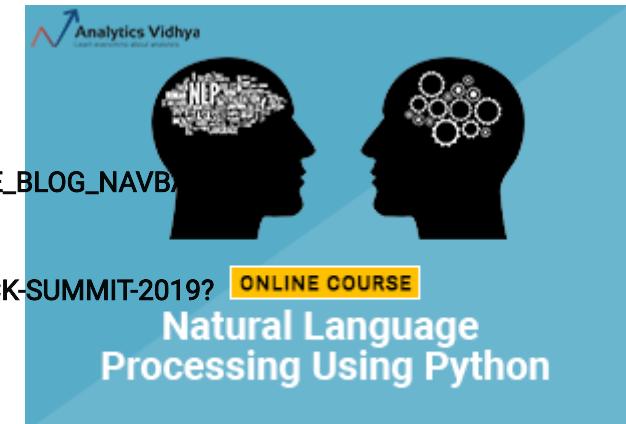
BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm\_source=HOME\_BLOG\_NAVBAR)
r):
COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM/) ~
    for action in range(environment.nA):
        for probability, next_state, reward, terminate
HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)
d in environment.P[state][action]:
    DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm\_source=HOME\_BLOG\_NAVBAR)
        (reward + discount_factor * V[next_state])
return action_values
DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm\_source=HOME\_BLOG\_NAVBAR)

```

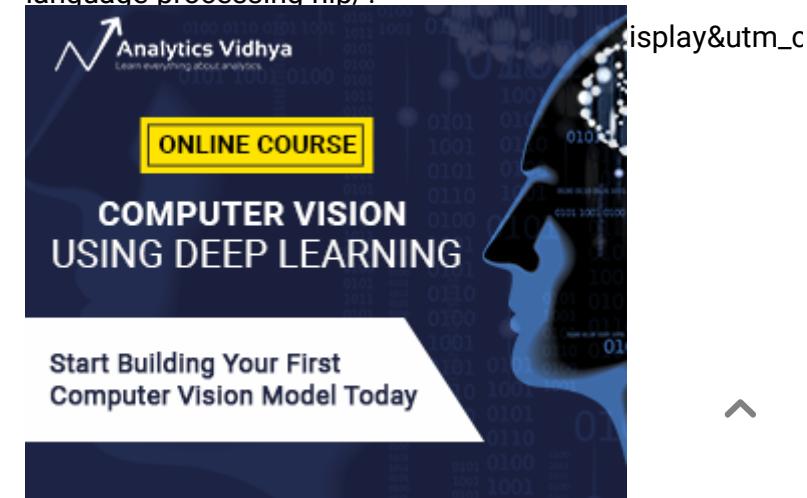
Now, the overall policy iteration would be as described below. This will return a tuple (policy,V) which is the optimal policy matrix and value function for each state.

CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))

([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/COMPUTER-VISION-USING-DEEP-LEARNING-VERSION2/?utm_source=blog&utm_medium=Stickybanner2](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2))



([HTTPS://COURSES.ANALYTICSVIDHYA.COM/COURSES/NATURAL-LANGUAGE PROCESSING-NLP/](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/))

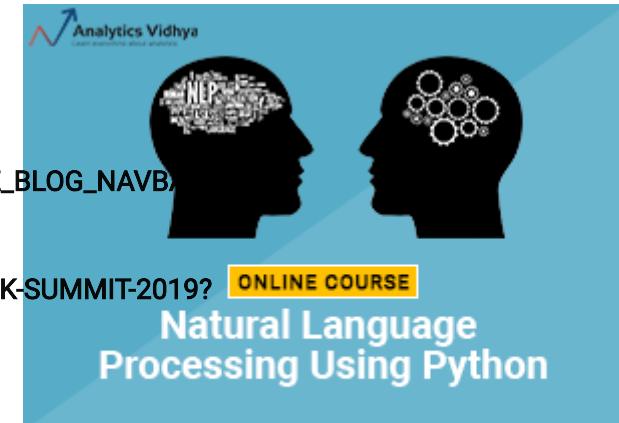


```

BLOG (https://www.analyticsvidhya.com/blog/?utm\_source=HOME\_BLOG\_NAVBAR)
rations=1e9):
COURSES#(https://courses.analyticsvidhya.com) ~
    #num states x num actions / num actions
    policy = np.ones((environment.nS, environment.nA)) / environment.nA
HACKATHONS (https://datahack.analyticsvidhya.com/contest/all)
DATAMIN# Initialize counter of evaluated policies
evaluated_policies = 1
# Repeat until convergence or critical number of iterations reached
DATAHACK SUMMIT 2019 (https://www.analyticsvidhya.com/datahack-summit-2019?utm\_source=HOME\_BLOG\_NAVBAR)
CONTACT (https://www.analyticsvidhya.com/contact/)
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

```

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=display&utm_medium=Stickybanner2)



```

# We will try every possible action in
BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm_source=HOME_BLOG_NAVBAR)
a current state

action_value = one_step_lookahead(environment, state, V, discount_factor)

# Select a better action
COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ^
environment, state, V, discount_factor)

HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)
best_action = np.argmax(action_value)

# If action didn't change
DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR)
if current_action != best_action:
    stable_policy = True

# Greedy policy update
DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR)
policy[state] = np.eye(environment.n_actions)[best_action]

evaluated_policies += 1

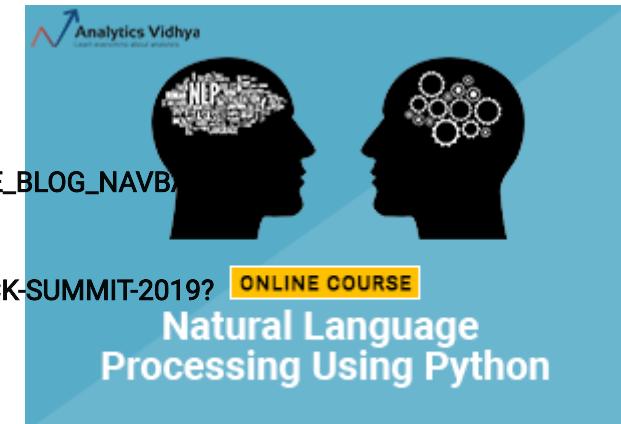
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/) is not
changing anymore, then return final policy and value function

if stable_policy:
    print(f'Evaluated {evaluated_policies} policies.')
    return policy, V

```

Value Iteration in python

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=HOME_BLOG_NAVBAR)



Start Building Your First
Computer Vision Model Today

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

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ▾

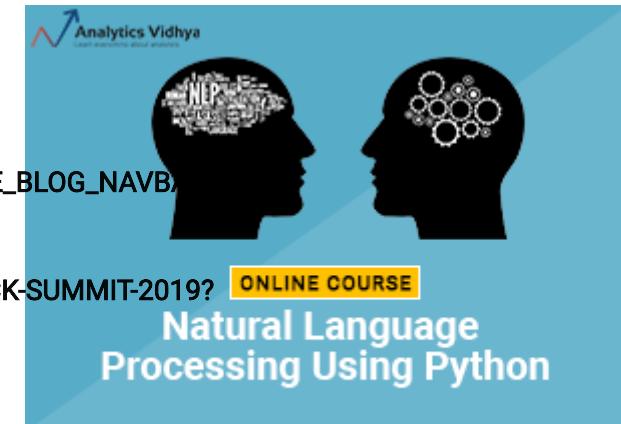
HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR](https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar))

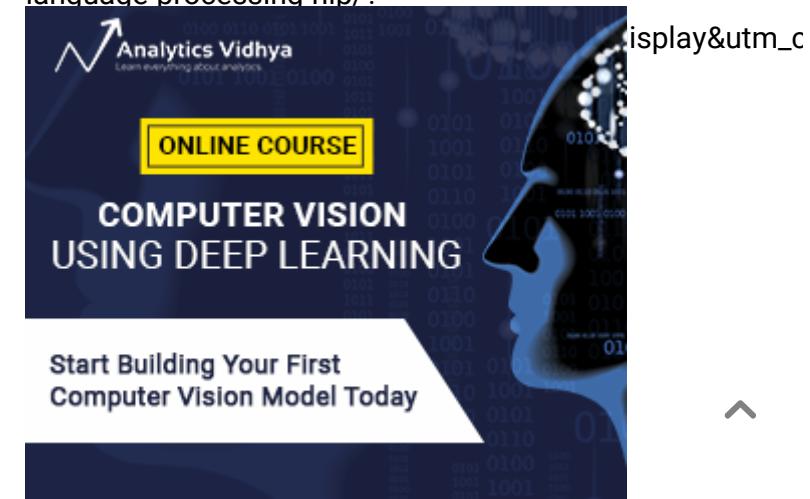
DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=home_blog_navbar))

UTM_SOURCE=HOME_BLOG_NAVBAR)

CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))



(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



```

BLOG (https://www.analyticsvidhya.com/blog/?utm\_source=HOME\_BLOG\_NAVBAR)
-9, max_iterations=1e9):

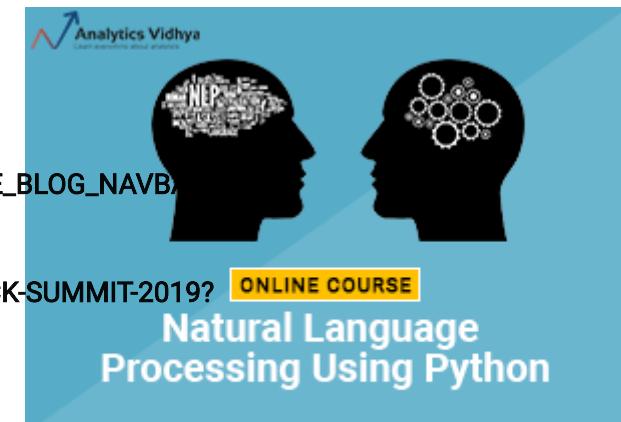
COURSES (https://courses.analyticsvidhya.com)
# Initialize state-action function with zeros for each
environment state

V = np.zeros(environment.nS)

HACKATHONS (https://datahack.analyticsvidhya.com/contest/all)
for i in range(int(max_iterations)):
    DATAMIN (# Early stopping condition
              https://datamin.analyticsvidhya.com/?utm\_source=HOME\_BLOG\_NAVBAR)
              delta = 0
    # Update each state
    DATAHACK SUMMIT 2019 (https://www.analyticsvidhya.com/datahack-summit-2019?utm\_source=HOME\_BLOG\_NAVBAR)
    for state in range(environment.nS):
        # Do a one-step lookahead to calculate
        # state-action values
        UTM_SOURCE=HOME_BLOG_NAVBAR
        action_value = one_step_lookahead(environment, state, V, discount_factor)
CONTACT (https://www.analyticsvidhya.com/contact/)
        # Select best action to perform based
        on the highest state-action value
        best_action_value = np.max(action_value)
    e)
        # Calculate change in value
        delta = max(delta, np.abs(V[state] - best_action_value))
        # Update the value function for current
        t state
        V[state] = best_action_value

```

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=display&utm_medium=Stickybanner2)



```

# Check if we can stop
BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm_source=HOME_BLOG_NAVBAR)
    if delta < theta:
        print(f'Value-iteration converged at i')
COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ~
iteration#{i}.')
        break
HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)

# Create a deterministic policy using the optimal value
DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR)
    function
        policy = np.zeros([environment.nS, environment.nA])
DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=blog&utm_medium=Stickybanner2/)

# One step lookahead to find the best action for
# this state
UTM_SOURCE=HOME_BLOG_NAVBAR)

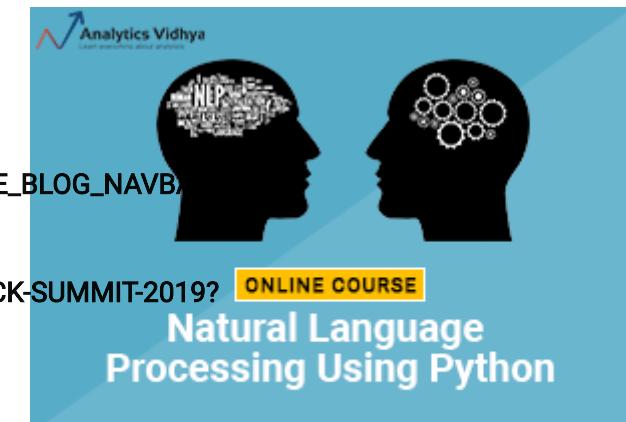
action_value = one_step_lookahead(environment,
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

# 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

```

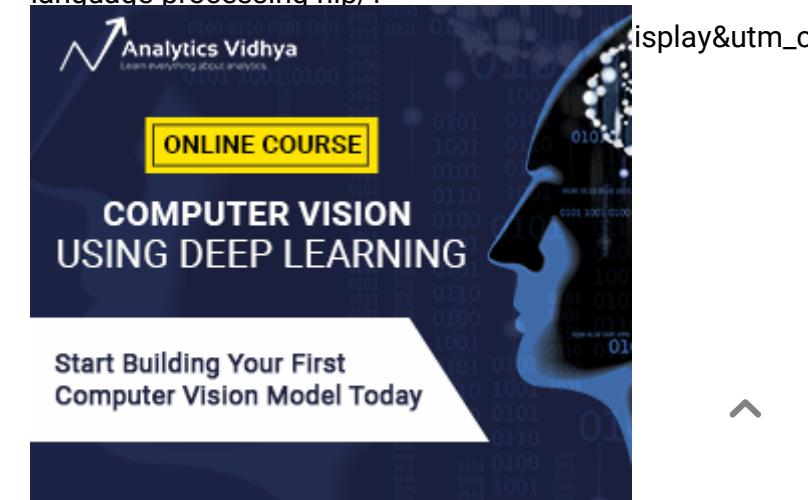
Finally, let's compare both methods to look at which of them works better in a practical setting. To do this, we will try to learn the optimal policy for the frozen lake environment using both techniques

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



described above. Later, we will check which technique performed better based on the average return after 10,000 episodes.

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>?utm_source=blog&utm_medium=Stickybanner2)

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ▾

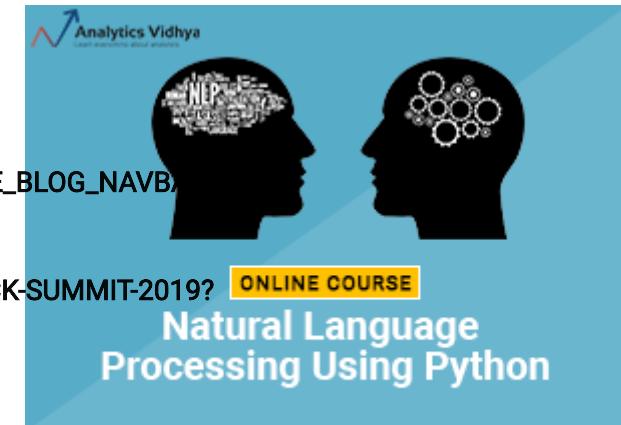
HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

DATAMIN ([HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=home_blog_navbar](https://datamin.analyticsvidhya.com/?utm_source=home_blog_navbar))

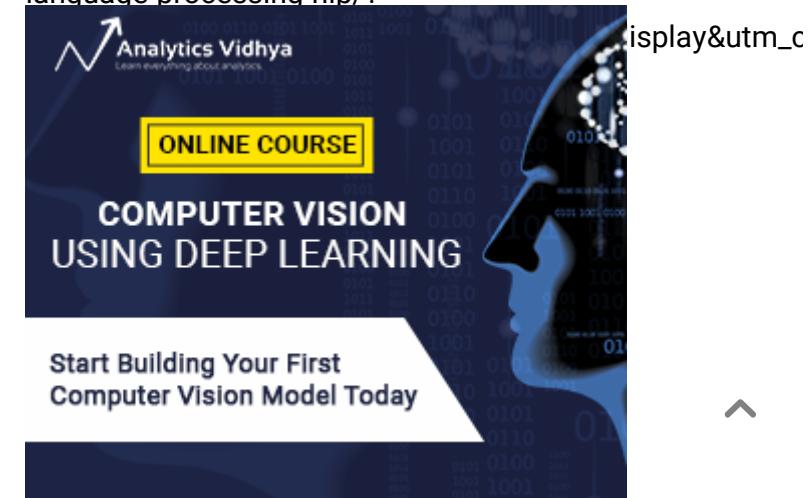
DATAHACK SUMMIT 2019 ([HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=home_blog_navbar](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=home_blog_navbar))

UTM_SOURCE=HOME_BLOG_NAVBAR)

CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))



Learn to Solve
Text Classification Problems Using **NLP**
(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>?)



```

BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm_source=HOME_BLOG_NAVBAR)
wins = 0

COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ~
    for episode in range(n_episodes):
        terminated = False

HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)
    state = environment.reset()

DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR)
    # Select best action to perform in a c

urrent state
DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR)
    action = np.argmax(policy[state])

    # Perform an action and observe how env
    environment.step(action)

    next_state, reward, terminated, info =
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)
environment.step(action)

    # Summarize total reward
    total_reward += reward

    # Update current state
    state = next_state

    # Calculate number of wins over episod
es

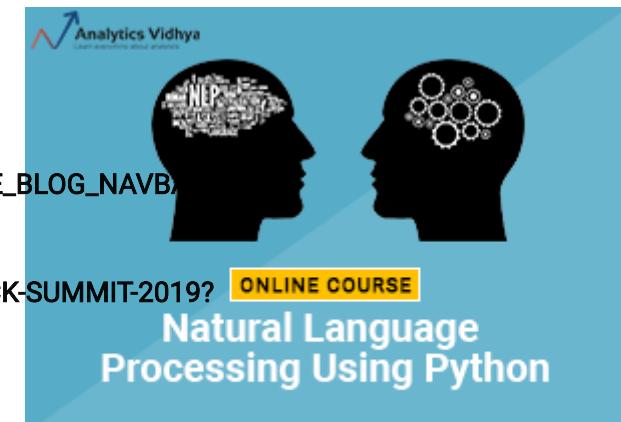
    if terminated and reward == 1.0:
        wins += 1

average_reward = total_reward / n_episodes

return wins, total_reward, average_reward

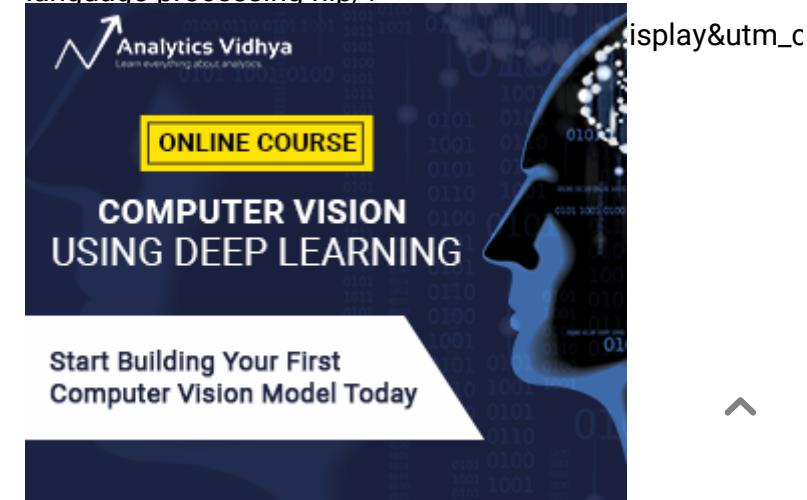
```

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>?utm_medium=Stickybanner2)



```

BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm_source=HOME_BLOG_NAVBAR)
# Number of episodes to play
n_episodes = 10000
COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM) ^
# Functions to find best policy
solvers = [('Policy Iteration', policy_iteration),
HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)
    ('Value Iteration', value_iteration)]
for iteration_name, iteration_func in solvers:
DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR)
    # Load a Frozen Lake Environment
    environment = gym.make('FrozenLake-v0')
DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR)
    # Search for an optimal policy using policy iteration
    policy, V = iteration_func(environment.env)
UTM_SOURCE=HOME_BLOG_NAVBAR
    # Apply best policy to the real environment
    wins, total_reward, average_reward = play_episodes(environment,
CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)
        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')

```

```

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

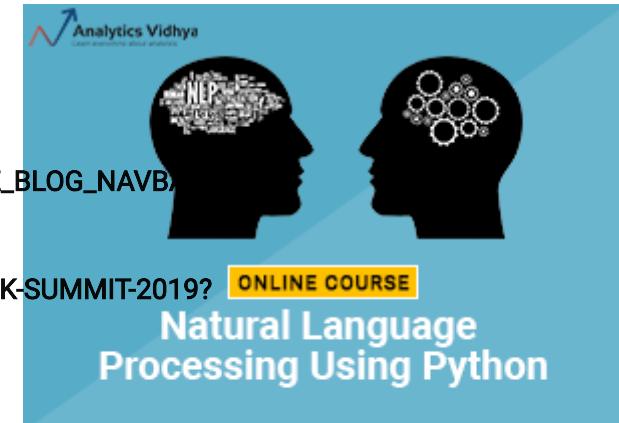
```

```

Value-iteration converged at iteration#523.
Value Iteration :: number of wins over 10000 episodes = 7397
Value Iteration :: average reward over 10000 episodes = 0.7397

```

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using **NLP**

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>?utm_source=blog&utm_medium=Stickybanner2)



Start Building Your First
Computer Vision Model Today

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

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

COURSES ([HTTPS://COURSES.ANALYTICSVIDHYA.COM](https://courses.analyticsvidhya.com)) ▾

End Notes

HACKATHONS ([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL](https://datahack.analyticsvidhya.com/contest/all))

In this article, we became familiar with model based planning using [DATAHACK SUMMIT 2019 \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/SUMMIT?utm_source=HOME_BLOG_NAVBAR\)](https://datahack.analyticsvidhya.com/summit/2019?utm_source=home_blog_navbar). If you, specifically in this 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 (https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2) covering different algorithms within this exciting domain.

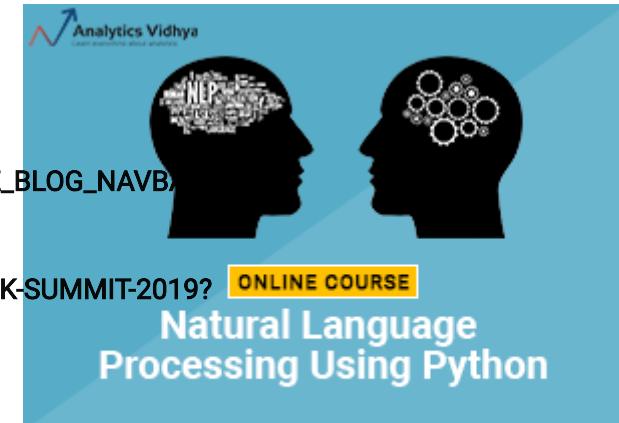
CONTACT ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))

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_other_global_all_co_prtnr_py_PartBadge-Mar2515-1

Share this:



Learn to Solve
Text Classification Problems Using NLP

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blog&utm_medium=Stickybanner2)



ONLINE COURSE
COMPUTER VISION
USING DEEP LEARNING

Start Building Your First
Computer Vision Model Today

[!\[\]\(0f4ae20874db623ff48a27215649ce03_img.jpg\) \(https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=HOME_BLOG_NAVBAR)

[!\[\]\(f4a0621425c41ac5740861def630d6a7_img.jpg\) \(https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=COURSES\)](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=COURSES)

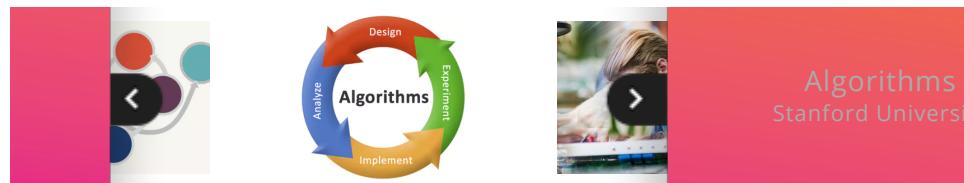
[!\[\]\(9b4000408f3699de3a705e848e6ae947_img.jpg\) \(https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=HACKATHONS\)](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=HACKATHONS)

[!\[\]\(beadafdc0beb7d8dd0a09f518e768281_img.jpg\) \(https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=DATAMIN\)](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=DATAMIN)

[!\[\]\(3f41268aaa93dab4a01c59fb2f124f87_img.jpg\) \(https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=DATAHACK_SUMMIT\)](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/?utm_source=DATAHACK_SUMMIT)

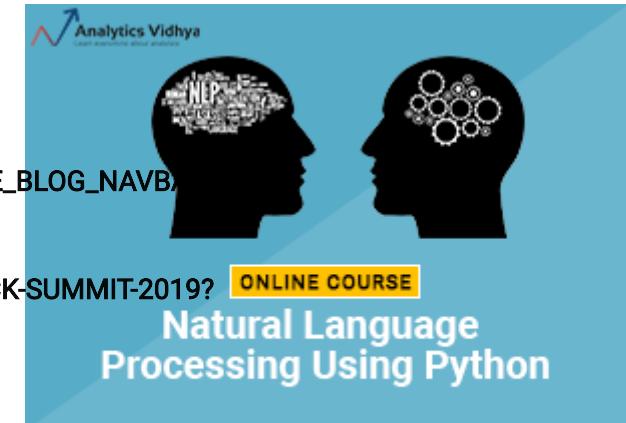
UTM_SOURCE=HOME_BLOG_NAVBAR
Like this:

Loading...
[!\[\]\(34543dd4ff7f078317aba2ea094681a5_img.jpg\) CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact)

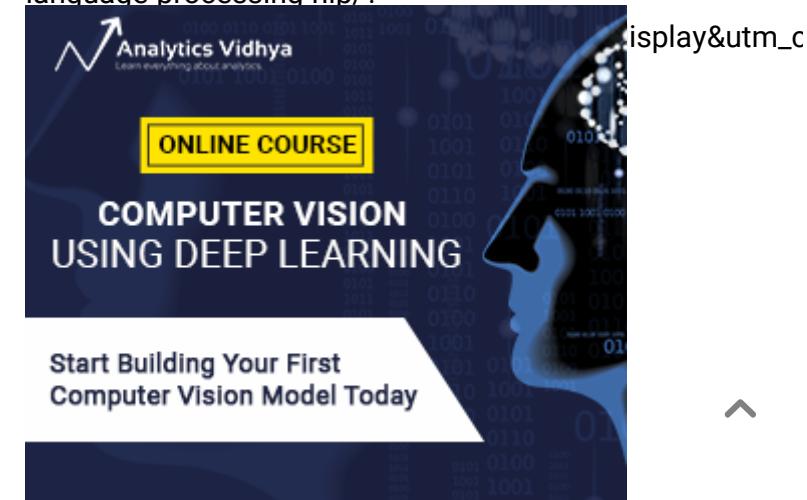


TAGS : [PYTHON \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/PYTHON/\)](https://www.analyticsvidhya.com/blog/tag/python/),
[REINFORCEMENT LEARNING
\(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/REINFORCEMENT-LEARNING/\)](https://www.analyticsvidhya.com/blog/tag/reinforcement-learning/)

[!\[\]\(9cbf1bb4a206d9681b6d411f6f46a945_img.jpg\) \(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2\)](https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve Text Classification Problems Using NLP
[\(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)



BLOG (<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/>)

utm_source=HOME_BLOG_NAVBAR

◀ DataHack Radio ... Performing Speech ▶

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

COURSES (<https://courses.analyticsvidhya.com/>)

and Object Recognition using just One Model System

HACKATHONS WITH MIT'S ML CONTEST/ALL

DATAMIN (<https://datamin.analyticsvidhya.com/>)

(<https://www.analyticsvidhya.com/blog/2018/09/dataladhyva.com/blog/radio-data-science-podcast-jeanette-wing/>)

DATAHACK SUMMIT 2019 (<https://www.analyticsvidhya.com/datahack-summit-2019/>)

UTM_SOURCE=HOME_BLOG_NAVBAR



CONTACT (<https://www.analyticsvidhya.com/contact/>)

Ankit Choudhary

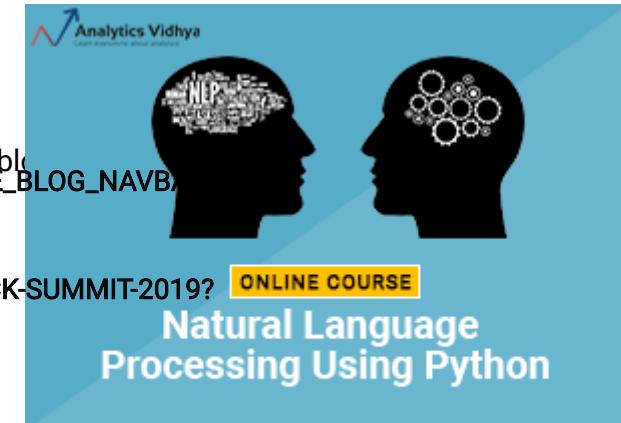
(<https://www.analyticsvidhya.com/blog/author/ankit2106/>).

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/>).

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>)

utm_source=blog&utm_medium=Stickybanner2)



Learn to Solve
Text Classification Problems Using NLP

(<https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



Start Building Your First
Computer Vision Model Today

[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR)

(https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/?utm_source=blog&utm_medium=Stickybanner2)

[COURSES \(HTTPS://COURSES.ANALYTICSVIDHYA.COM\)](https://courses.analyticsvidhya.com/) ▾

This article is quite old and you might not get a prompt response from the author. We request you to post this comment on [HACKATHONS \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL\)](https://datahack.analyticsvidhya.com/contest/all) Analytics Vidhya's [Discussion portal](#)

(<https://discuss.analyticsvidhya.com/>) to get your queries

[DATAMIN \(HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://datamin.analyticsvidhya.com/?utm_source=HOME_BLOG_NAVBAR)

4 COMMENTS

[DATAHACK SUMMIT 2019 \(HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/datahack-summit-2019?utm_source=HOME_BLOG_NAVBAR)

 **RAMESH MATHIKUMAR**
[utm_source=HOME_BLOG_NAVBAR](https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/#comment-155008)

[September 18, 2018 at 11:11 pm](#)

[Reply](#)

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

Good Article Ankit.



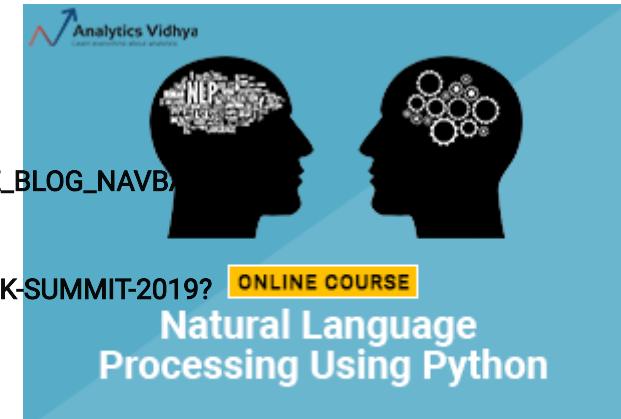
BHUMIKA

[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?



Learn to Solve
Text Classification Problems Using **NLP**

(https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blog&utm_medium=Stickybanner2)





ANKIT CHOUDHARY
[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?UTM_SOURCE=HOME_BLOG_NAVBAR\)](https://www.analyticsvidhya.com/blog/?utm_source=HOME_BLOG_NAVBAR)
September 24, 2018 at 12:20 pm

[Reply](#)

(<https://courses.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/>?utm_source=blog&utm_medium=Stickybanner2)

[COURSES \(HTTPS://COURSES.ANALYTICSVIDHYA.COM\)](https://www.analyticsvidhya.com/courses/computer-vision-using-deep-learning-version2/) ▾
 (Solving Model Based Planning Example)
[programming/#comment-155060](#))

HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)
 You can refer to this stack overflow query:

<https://stats.stackexchange.com/questions/243384/deriving-the-bellmans-equation-in-reinforcement-learning>

[CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/?utm_source=HOME_BLOG_NAVBAR)

(<https://www.analyticsvidhya.com/blog/2018/09/reinforcement-learning-model-based-planning-dynamic-programming/#comment-157693>)

[UTM_SOURCE=HOME_BLOG_NAVBAR](#)

VIJIT

[Reply](#)

[CONTACT \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/)
 (<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.

Natural Language Processing Using Python

Learn to Solve Text Classification Problems Using **NLP**

COMPUTER VISION USING DEEP LEARNING

Start Building Your First Computer Vision Model Today

<p>BLOG (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/?utm_source=HOME_BLOG_NAVBAR)</p> <p>ANALYTICS VIDHYA</p> <p>About Us (http://www.analyticsvidhya.com/about-me/)</p> <p>Our Team (http://www.analyticsvidhya.com/team/)</p> <p>Career (http://www.analyticsvidhya.com/career/)</p> <p>Contact Us (http://www.analyticsvidhya.com/contact/)</p> <p>COURSES (HTTPS://COURSES.ANALYTICSVIDHYA.COM/)</p> <p>HACKATHONS (HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/ALL)</p> <p>DATAMIN (HTTPS://DATAMIN.ANALYTICSVIDHYA.COM/?utm_source=HOME_BLOG_NAVBAR)</p> <p>DATAHACK SUMMIT 2019 (HTTPS://WWW.ANALYTICSVIDHYA.COM/DATAHACK-SUMMIT-2019?utm_source=HOME_BLOG_NAVBAR)</p> <p>CONTACT (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)</p>	<p>DATA SCIENTISTS</p> <p>Blog (https://www.analyticsvidhya.com/blog/)</p> <p>Hackathon (https://datahack.analyticsvidhya.com/)</p> <p>Discussions (https://discuss.analyticsvidhya.com/)</p> <p>Apply Jobs (https://www.analyticsvidhya.com/jobs/)</p> <p>Leaderboard (https://datahack.analyticsvidhya.com/leaderboard/)</p> <p>COMPANIES</p> <p>Post Jobs (https://www.analyticsvidhya.com/corporate/)</p> <p>Trainings (https://trainings.analyticsvidhya.com/)</p> <p>Advertising (https://www.analyticsvidhya.com/advertising/)</p> <p>Reach Us (https://www.analyticsvidhya.com/contact/)</p> <p>JOIN OUR COMMUNITY</p> <p>f (https://www.facebook.com/analyticsvidhya) 22062</p> <p>t (https://twitter.com/analyticsvidhya)</p> <p>g+ (https://plus.google.com/+Analyticsvidhya) Followers</p> <p>l (https://in.linkedin.com/in/analyticsvidhya) Followers</p> <p>n (https://plus.google.com/m/+Analyticsvidhya) Followers</p> <p>Natural Language Processing Using Python</p> <p>Followers</p> <p>Learn to Solve Text Classification Problems Using NLP</p> <p>COMPUTER VISION USING DEEP LEARNING</p> <p>Start Building Your First Computer Vision Model Today</p>
--	---

© Copyright 2013-2019 Analytics Vidhya.

[Privacy Policy](https://www.analyticsvidhya.com/privacy-policy/) (<https://www.analyticsvidhya.com/privacy-policy/>)

[Terms of Use](https://www.analyticsvidhya.com/terms-of-use/) (<https://www.analyticsvidhya.com/terms-of-use/>)

[Refund Policy](https://www.analyticsvidhya.com/refund-policy/) (<https://www.analyticsvidhya.com/refund-policy/>)

Don't have an account? Sign up (<https://id.analyticsvidhya.com/accounts/signup/>) here