temp state=deepcopy(self.state) if(action==0): if(self.state[1]>0): temp\_state[1]-=1 if (action==1): if(self.state[1] < self.size-1):</pre> temp\_state[1]+=1 if (action==2): if(self.state[0] < self.size-1):</pre> temp state[0]+=1if(action==3): if(self.state[0]>0): temp\_state[0]-=1 if(update==1): self.state=temp state if(temp\_state==self.ter): return temp state, self.ter reward, 1 return temp state, self.non reward, 0 def reset(self): self.state=self.start self.grid=[[0 for i in range(self.size)] for j in range(self.size)] In [3]: class DiscreteAgent(): def \_\_init\_\_(self,env): self.env=env self.sweep no=0 self.policy=[[1 for i in range(env.size)] for j in range(env.size)] self.another\_policy=[0 for i in range(env.size)] def update(self): pass def get\_action(self, state): In [4]: class ValueIteration(DiscreteAgent): def update(self): previousreward = deepcopy(self.env.grid) is\_converged=**False** for i in range(0, self.env.size): for j in range(0, self.env.size): values=[] for action in self.env.actions: print("krishna") self.env.state=[i,j] for action1 in self.env.actions: self.env.state=[i,j] state,r,x=self.env.step(action1,0) print(str(action)+"->"+str(action1)+"->"+str(self.en v.state) + "->" + str(state) + "->" + str(r) + "->" + str(self.env.P[action][action1]))print(self.env.P[action][action1]) ans+=self.env.P[action][action1]\*(r+previousreward[sta te[0]][state[1]]\*self.env.gamma) values.append(r+previousreward[state[0]][state[1]]\*s elf.env.gamma) values.append(ans) print(values) if([i,j]!=self.env.ter): self.env.grid[i][j]=max(values) self.policy[i][j]=values.index(max(values)) error=max(error, abs(previousreward[i][j]-self.env.grid[i][j])) **if** (error<0.01): is converged=**True** self.sweep no+=1 return self.sweep no,is converged def get action(self, state): return self.policy[state[0]][state[1]] In [5]: env=DiscreteEnvironment([0,0],[7,7]) # print(env.grid) agent=ValueIteration(env) # agent=PolicyIteration(env) # print(agent.sweep no) sweep no, is converged = 0, False max sweeps = 49returns=[0 for i in range(max sweeps+1)] while sweep no<max sweeps:</pre> sweep no,stop=agent.update() is terminated=False state=env.start while not is terminated: temp=state action = agent.policy[state[0]][state[1]] env.state=temp next state, reward, is terminated = env.step(action,0) if(next state==temp): print(sweep no) print("krishna") reward=-5 returns[sweep no]+=reward returns[sweep no]+=reward print(reward) state=next state print(sweep no) print(returns) 3.5, 3.5] In [6]: # agent.evaluate policy() # while not is converged: while sweep no<max sweeps:</pre> sweep no, is converged = agent.update() # print(env.grid) # print(sweep no) print(agent.policy) 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 1]] In [7]: class PolicyIteration(DiscreteAgent): def update(self): temp=deepcopy(self.policy) self.evaluate\_policy() self.update policy() self.sweep no+=1 if(self.policy==temp): return self.sweep no,True return self.sweep\_no,False def evaluate\_policy(self): is\_converged=False while(is converged==False): while(self.sweep\_no<20):</pre> error=0 previousreward = deepcopy(self.env.grid) for i in range(0, self.env.size): for j in range(0, self.env.size): self.env.state=[i,j] action=self.policy[i][j] for action1 in self.env.actions: self.env.state=[i,j] state, r, x=self.env.step(action1,0) val+=self.env.P[action][action1]\*(r+previousreward[sta te[0]][state[1]]\*self.env.gamma) if([i,j]!=self.env.ter): self.env.grid[i][j]=val error=max(error,abs(previousreward[i][j]-self.env.grid[i][ j])) print (error) **if**(error<0.001): is\_converged=**True** def update\_policy(self): for i in range(0, self.env.size): for j in range(0, self.env.size): values=[] for action in self.env.actions: val=0 for action1 in self.env.actions: self.env.state=[i,j] state, r, x=self.env.step(action1,0) val+=self.env.P[action][action1]\*(r+self.env.grid[stat e[0]][state[1]]\*self.env.gamma) values.append(val) self.policy[i][j]=values.index(max(values)) def get action(self, state): return self.policy[state[0]][state[1]] In [8]: env1=DiscreteEnvironment([0,0],[7,7]) agent1=PolicyIteration(env1) i=0 stop=Falsesweep no=0 # max sprint(hile(stop=False): returns1=[0 for i in range(max sweeps+1)] while sweep no<max sweeps:</pre> sweep no,stop=agent1.update() is terminated=False state=env1.start while not is terminated: temp=state action = agent1.policy[state[0]][state[1]] env1.state=temp next\_state, reward, is\_terminated = env1.step(action,0) if(next state==temp): print(sweep no) print("krishna") reward=-5 returns1[sweep no]+=reward returns1[sweep no]+=reward print(reward) state=next state print(sweep\_no) print(returns1) 5, 3.5, 3.5, 3.5] In [9]: print(agent1.policy) # print(env1.grid) c=[[ 0 for i in range(agent.env.size)]for i in range(agent.env.size)] for i in range(8): for j in range(8): c[i][j]=agent.policy[i][j]-agent1.policy[i][j] 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 2], [1, 1, 1, 1, 1, 1, 1, 1]] [[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, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0]] In [10]: returns2=[0 for i in range(50)] for i in range (50): is terminated=False state=env1.start while not is terminated: temp=state action = random.randint(0,3) env.state=temp next\_state, reward, is\_terminated = env.step(action,0) print(str(state)+"->"+str(action)+"->"+str(next\_state)) returns2[i]+=reward state=next state print(returns2) [-192.5, -503.5, -24.0, -176.0, -84.0, -16.5, -506.5, -45.5, -18.0, -31.0,-4.0, -27.0, -440.5, -5.0, -34.0, -29.5, -272.0, -224.5, -21.5, -32.5, -23. 0, -50.0, -75.5, -76.5, -260.0, -205.0, -218.5, -132.5, -243.5, -12.5, -58 8.0, -87.5, -67.5, -455.0, -93.5, -297.5, -26.0, -203.5, -15.0, -418.5, -43.5, -311.5, -459.5, -536.5, -431.5, -54.0, -9.0, -53.0, -16.5, -74.0] In [11]: import matplotlib.pyplot as plt import matplotlib.patches as mpatches from matplotlib.colors import colorConverter as cc import numpy as np def plot\_mean\_and\_CI(mean, lb, ub, color\_mean=None, color\_shading=None): # plot the shaded range of the confidence intervals plt.fill\_between(range(mean.shape[0]), ub, lb, color=color\_shading, alpha=.5) # plot the mean on top plt.plot(mean, color mean) returns=np.array(returns) returns1=np.array(returns) # generate 3 sets of random means and confidence intervals to plot ub0 = returns + np.random.random(50) + .51b0 = returns - np.random.random(50) - .5ub1 = returns1 + np.random.random(50) + .5lb1 = returns1 - np.random.random(50) - .5returns3 =np.array(returns2)/100 #Note: Random Agents returns scaled by 100 for better visualisation ub2 = returns3 + np.random.random(50) + .51b2 = returns3 - np.random.random(50) - .5# plot the data fig = plt.figure(1, figsize=(7, 2.5)) plot\_mean\_and\_CI(returns, ub0, 1b0, color\_mean='k', color\_shading='k') plot\_mean\_and\_CI(returns1, ub1, lb1, color\_mean='b', color\_shading='b') plot\_mean\_and\_CI(returns3, ub2, lb2, color\_mean='g--', color\_shading='g') 5.0 2.5 -7.520 **States**  $\{i,j\}$  where i=[0,20], j=[0,20]0 is dummy state introduced for conveinience where i is number of cars in location1 and j is number of cars in locat ion2 **Actions** actions in state A=[-5,5]Here maximum number of cars that can be moved from one location to othe Instead of having two numbers to know from which location we moved to w hich location, We maintain a single number, negative quantity represents moving from lo cation2 to location1 and vice versa **Rewards** r(s',a,s) = min(ren1,s[0]-a) + min(ren2,s[1]+a) \*cwhere ren1=no of cars given for rent at location1 ren2=no of cars given for rent at location2 s[0] is number of cars at location1 s[1] is number of cars at location2 **Transitions**  $p({s[0]-a-ren1+ret1,s[0]+a-ren2+ret2}/s,a)=poisson(ren1,f req)*poisson$ (ren2,s req)\*poisson(ret1, f ret)\*poisson(ret2, s ret) ren1=no of cars given for rent at location1=min(asked for rent,availabl e after taking action a) ren2=no of cars given for rent at location2=min(asked for rent,availabl e after taking action a) ret1=no of cars returned at location1 ret2=no of cars returned at location2 poisson (ren1, f req) is probability of poission distribution for value r en1 where lambda=f req poisson(ren1, f req) is probability of poission distribution for value r en2 where lambda=s req poisson(ren1, f req) is probability of poission distribution for value ret1 where lambda=f ret poisson (ren1, f req) is probability of poission distribution for value r et2 where lambda=s ret and f\_ren,s\_ren,f\_ret,s\_ret are given **Bellman update equation** let c1=a+ren1-ret1 c2=a-ren2+ret2  $P=p({s[0]-c1,s[1]+c2}/s,a)$ Policy evaluation update: where a=poilcy(s)  $v(s)=\{P*(r(s',a,s)+v(s'))\}$ Policy update equation:  $policy(s) = (a) for \max_{orall a \in A} \{P*(r(s',a,s) + v(s'))\}$ In [108]: import math class CarEnvironment(): def \_\_init\_\_(self): self.maxcars=20 self.maxmove=5 self.f req=3 self.s req=4 self.f\_ret=3 self.s ret=2 self.rent=10 self.cost=-2self.up=9 self.prob=dict() self.state=[0,0] self.grid=[[0 for i in range(self.maxcars+1)] for j in range(self.maxc ars+1)] self.size=self.maxcars+1 self.actions=[i for i in range(-self.maxmove, self.maxmove+1)] self.gamma=0.9 #precalculating and storing in dictionary as this is required so many numb er of times and will become time consuming def poisson(self,x, 1): key = x \* 10 + 1if key not in self.prob.keys(): self.prob[key] = np.exp(-1) \* pow(1, x) / math.factorial(x)return self.prob[key] In [109]: class PolicyIteration2(DiscreteAgent): def update(self): temp=deepcopy(self.policy) print("krishna1") self.evaluate\_policy() print("krishna2") self.update\_policy() print("krishna3") self.sweep\_no+=1 if (self.policy==temp): return self.sweep no,True print("krishna4") return self.sweep\_no,False def evaluate\_policy(self): is converged = **False** #no of iterations to run for evaluating policy while(sweep<30):</pre> print("kmt1") error=0 #storing current environment state prevstate=deepcopy(self.env.grid) for i in range(self.env.size): for j in range(self.env.size): print("kmt2") self.env.state=[i,j] action=self.policy[i][j] #max cars at that location cant exceed 20 cars1 = int(min(self.env.state[0] - action, self.env.maxca rs)) cars2 = int(min(self.env.state[1] + action, self.env.maxca rs)) #checking for all possibilities of number of cars that cou ld be rented(as poisson has all probabilities to infinity neglecting probabilit ies after 9) for ren1 in range(0, self.env.up): for ren2 in range(0, self.env.up): print("kmt3") rentprob = self.env.poisson(ren1, self.env.f req) \* self.env.poisson(ren2, self.env.s req) totalren1 = min(cars1, ren1) totalren2 = min(cars2, ren2) rew = ((totalren1 + totalren2) \* self.env.rent+sel f.env.cost\* abs(action))\*rentprob #checking for all possibilities of no of cars returned as to know the next state for ret1 in range(0, self.env.up): for ret2 in range(0, self.env.up): print("kmt4") retprob = self.env.poisson(ret1, self.env. f ret) \* self.env.poisson(ret2,self.env.s req) prob=retprob\*rentprob cars1 end = min(cars1 - totalren1 + ret1, self.env.maxcars) cars2\_end = min(cars2 - totalren2 + ret2, self.env.maxcars) #rew is r(s,a) and returns is total value of state returns += retprob \* rew + prob\*(self.env. gamma\* prevstate[cars1 end][cars2 end]) self.env.grid[i][j]=returns error=max(error, abs(prevstate[i][j]-self.env.grid[ i][j])) sweep+=1 **if** (error<0.01): is converged=True def update policy(self): for i in range(0, self.env.size): for j in range(0, self.env.size): vals = [] $\hbox{\it\#for all states checking which action gives maximum possible $v$}$ alue for action in self.env.actions: #taking only for possible actions because you cannot move more than nof of cars at that state if ((action >= 0 and i >= action) or (action < 0 and j >= abs(action))): self.env.state=[i,j] returns=0 cars1 = int(min(self.env.state[0] - action, self.env.m axcars)) cars2 = int(min(self.env.state[1] + action, self.env.m axcars)) for ren1 in range(0, self.env.up): for ren2 in range(0, self.env.up): rentprob = self.env.poisson(ren1, self.env.f\_r eq) \* self.env.poisson(ren2, self.env.s\_req) totalren1 = min(cars1, ren1) totalren2 = min(cars2, ren2) rew = ((totalren1 + totalren2) \* self.env.rent +self.env.cost\* abs(action))\*rentprob for ret1 in range(0, self.env.up): for ret2 in range(0, self.env.up): retprob = self.env.poisson(ret1, self. env.f ret) \* self.env.poisson(ret2, self.env.s req) prob=retprob\*rentprob cars1\_end = min(cars1 - totalren1 + re t1, self.env.maxcars) cars2\_end = min(cars2 - totalren2 + re t2, self.env.maxcars) returns += retprob \* rew + prob\*(self. env.gamma\* self.env.grid[cars1 end][cars2 end]) vals.append(returns) bestAction = vals.index(max(vals)) self.policy[i] [j] = self.env.actions[bestAction] def get\_action(self, state): return self.policy[state[0]][state[1]] In [110]: env4=CarEnvironment() agent4=PolicyIteration2(env4)  $\verb|is_converged=False|\\$ sweep no=0 max sweeps=1 while not is\_converged: sweep\_no,is\_converged=agent4.update() krishna1 krishna2 krishna3 krishna4 krishna1 krishna2 krishna3 krishna4 krishna1 krishna2 krishna3 krishna4 krishna1 krishna2 krishna3 In [112]: # print(env4.grid) # print(agent4.policy) # import numpy from mpl\_toolkits import mplot3d import numpy as np import matplotlib.pyplot as plt fig = plt.figure() ax = plt.axes(projection="3d") z points = []  $x_points = []$ y points = [] for i in range(21): for j in range (21): x points.append(i) y\_points.append(j) z\_points.append(env4.grid[i][j]) ax.scatter3D(x\_points, y\_points, z\_points, c=z\_points); #similarly plot for policy also 420 400 380 360  $^{0.0}_{\ 2.5}_{\ 5.0}_{\ 7.5}{}_{10.0}_{12.5}{}_{15.0}{}_{17.5}{}_{20.0}$ **States** but as you told i added two dummy states for convinience 0 and 100 Actions actions in state A  $s=\{0,1,..min(s,100-s)\}$ Rewards r(s',a,s)=1 if s'=100=0 otherwise **Transitions** p(s+a/s,a) = php(s-a/s,a)=1-phwhere ph= probability of getting heads on a toss **Bellman update equation**  $v(s) = \max \left\{ ph * (reward[s+a] + v(s+a)) + (1-ph) * (reward[s-a] + v(s-a)) \right\}$ In [118]: class Gamble\_environment(): def init (self,ph): self.grid=[0 for i in range(101)] self.state=1 self.p=ph self.ter=100 self.size=101 self.actions=[i for i in range(101)] self.rewards=[0 for i in range(101)] self.rewards[100]=1 self.gamma=1 In [119]: class ValueIteration1(DiscreteAgent): def update(self): previous = deepcopy(self.env.grid) is converged=**False** error=0 for i in range(self.env.size): values=[] for action in self.env.actions: **if** (action **in** range  $(1, \min(i, 100-i)+1)$ ): self.env.state=i ans+=self.env.p\*(self.env.rewards[i+action]+previous[i+act ion]\*self.env.gamma)+(1-self.env.p)\*(self.env.rewards[i-action]+previous[i-act ion]\*self.env.gamma) values.append(ans) if(i!=self.env.ter and i!=0): self.env.grid[i]=max(values) error=max(error, abs(previous[i]-self.env.grid[i])) **if**(error<0.0001): is converged=**True** self.sweep no+=1 return self.sweep no, is converged def get action(self, state): for state in range(1,101): cur=self.env.grid p=self.env.p val=0 for action in self.env.actions: if (action <=min(state, 100-state) and action>=1): ans=(p\*(self.env.rewards[state+action]+cur[state+action])+(1-p ) \* (self.env.rewards[state-action]+cur[state-action])) if(ans>val): val=ans r=action return r

In [120]: env2=Gamble environment (0.25) #change value to 0.15 and 0.65 for their plots

[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2,

# print(env.grid)

max sweeps = 1000

# print(env2.grid)(

for i in range (1, 101):

policy=[]

print(policy)

while not is converged:

agent2=ValueIteration1(env2)

# while sweep no<max sweeps:</pre>

sweep no, is converged = 0, False

sweep no, is converged = agent2.update()

policy.append(agent2.get action(i))

agent2.another policy=agent2.get action(i)

In [1]: import numpy as np

import random

0.1,0.1,0.7]]

In [2]: class DiscreteEnvironment():

import matplotlib.pyplot as plt

def init (self, start, stop):

self.grid[7][7]=10
self.non\_reward=-0.5
self.ter reward=10

self.actions=[0,1,2,3]
def step(self,action,update):

#change probability structure

self.P=[[1,0,0,0],[0,1,0,0],[0,0,1,0],[0,0,0,1]]

self.actions=["up","down","right","left"]

self.grid=[[0 for i in range(self.size)] for j in range(self.size)]

self.P=[[0.7,0.1,0.1,0.1],[0.1,0.7,0.1,0.1],[0.1,0.1,0.7,0.1],[0.1,

from copy import deepcopy

self.size=8

self.ter=stop

self.gamma=0.9

self.start=start
self.state=self.start