	P[2,0,2]=0.4 $P[1,1,2]=0.9$ $P[2,1,2]=0.1$ $R[0,0]=-0.4$ $R[0,1]=0$ $R[1,0]=2$ $R[1,1]=0$ $R[2,0]=-1$ $R[2,1]=-0.5$ $P is the transition 3x2x3 matrix of the given MDP model, the first index indicates the next state, the second one indicates the$
	action taken and the third one is for the previous (or the current) state. R is the reward 3x2 matrix, its rows represent the states and its columns the actions. Q2: Implement and run value iteration (VI) in order to identify a 0.01–optimal policy.
In [351]:	<pre>def VI(P,R,discount_factor,precision): v_new=np.zeros((1,np.shape(P)[0])) v_old=v_new+1 iteration=0 while(abs(np.max(v_old-v_new))>precision): v_old=v_new V=np.zeros((np.shape(P)[0],np.shape(P)[1])) for x in range(np.shape(P)[0]): for a in range(np.shape(P)[1]):</pre>
In [94]:	The function VI (value iteration) takes in argument the transition and reward matrices, The discount_factor (gamma) and the precision of the optimal policy. It returns an array of three elements, the first one is the optimal policy, the second one the optimal value and the third one is the number of iteration till the while condition is unsatisfied anymore. print('The 0.01-optimal policy is : ', VI(P,R,0.95,0.01)[0], '\nThe 0.01-optimal value is:', VI(P,R,0.95,0.01)[1]) The 0.01-optimal policy is : [1. 0. 1.] The 0.01-optimal value is: [12.20945133 14.20945133 12.86174503]
In [354]:	Plot vk - v* ∞ as a function of iteration k, knowing that d* = [a1, a0, a1] (i.e., implement policy evaluation). pi_optimal=[1,0,1] P_pi=np.array([P[:,int(pi_optimal[0]),0],P[:,int(pi_optimal[1]),1],P[:,int(pi_optimal[2]),2]]) r_pi=np.array([R[0,int(pi_optimal[0])],R[1,int(pi_optimal[1])],R[2,int(pi_optimal[2])]]) v_optimal=np.dot(np.linalg.inv(np.eye(np.shape(P)[0])-0.95*P_pi),r_pi)
In [355]:	Calculating the optimal value given the optimal policy, using the exact formula (matrix formula: $V\pi = (I - \gamma P\pi) - 1r\pi$). $\begin{aligned} \text{def policy} &= \text{valuation}(k, P, R, \text{discount} &= \text{factor}, \text{voptimal}): \\ \text{vonew} &= \text{np.zeros}((1, \text{np.shape}(P)[0])) \\ \text{for i in } &= \text{range}(k): \\ \text{Vonep.zeros}((\text{np.shape}(P)[0], \text{np.shape}(P)[1])) \\ \text{for x in } &= \text{range}(\text{np.shape}(P)[0]): \end{aligned}$
	for a in range(np.shape(P)[1]): V[x,a]=R[x,a]+discount_factor*np.sum(P[:,a,x]*v_new) v_new=np.max(V,axis=1) return abs(np.max(v_optimal-v_new)) The function policy_evaluation takes in argument The number of iterations k, the transition and reward matrices, The discount_factor(gamma) and the the optimal value. It returns the maximum difference between the elementwise of the optimal value and the estimated one.
	<pre>k=[i for i in range(300)] value_convergence=[policy_evaluation(i,P,R,0.95,v_optimal) for i in range(300)] plt.plot(k,value_convergence) plt.xlabel('Iterations') plt.ylabel(' v-v* ') plt.grid() plt.title('Figure1: Gap between the optimal value and the calculated one')</pre>
In [98]:	Figure1: Gap between the optimal value and the calculated one 14 12 10 8 8
	The figure tells us that the gap between the optimal value and the calculated one decreases when the number of iteration
	Q3: implement policy iteration (PI) with arbitrary initial policy. def PI (P,R,discount factor):
	<pre>p=np.zeros(np.shape(P)[0]) v_old=np.zeros(len(P)) is_changed=True iteration=0 while(is_changed): P_pi=np.array([P[:,int(p[0]),0],P[:,int(p[1]),1],P[:,int(p[2]),2]]) r_pi=np.array([R[0,int(p[0])],R[1,int(p[1])],R[2,int(p[2])]]) v_pi_new=np.dot(np.linalg.inv(np.eye(np.shape(P)[0])-discount_factor*P_pi),r_pi) if((v_pi_new==v_old).all()): is_changed=False v_old=v_pi_new V=np.zeros((np.shape(P)[0],np.shape(P)[1]))</pre>
	<pre>for x in range(np.shape(P)[0]): for a in range(np.shape(P)[1]): V[x,a]=R[x,a]+discount_factor*np.sum(P[:,a,x]*v_pi_new) p=np.argmax(V,axis=1) iteration+=1 return [p,v_pi_new,iteration]</pre> The function PI (Policy iteration) takes in argument The number of iterations k, the transition and reward matrices, The discount_factor (gamma). It returns an array of three elements, the first one is the optimal policy, the second one the optimal
T. [100]	value and the third one is the number of iteration till the while condition is unsatisfied anymore. For the policy evaluation step I used the matrix closed form, since the given MDP is not too complex (matrix formula: Vπ =(I-γPπ)-1rπ). Compare the speed of convergence w.r.t. VI and discuss the relative merits of the two approaches.
In [109]: In [110]:	t=time.time() value_iteration=VI(P,R,0.95,0) elapsed=time.time()-t t=time.time() policy_iteration=PI(P,R,0.95) elapsed_2=time.time()-t print('The time taken by the value iteration algorithm to converge is', elapsed) print('The number of value iterations till convergence is: ', value_iteration[2]) print('The time taken by the policy iteration algorithm to converge is', elapsed_2) print('The number of policy iterations till convergence is: ', policy_iteration[2]) The time taken by the value iteration algorithm to converge is 0.03574800491333008 The number of value iterations till convergence is: 642 The time taken by the policy iteration algorithm to converge is 0.0011980533599853516 The number of policy iterations till convergence is: 3
	The time taken to converge, in this case, of the VI approach is larger than the one taken by PI approach, but this is because the VI is an approximatly method and then if we are searching for exact convergences we should wait a lot of iterations unlike the PI approach which use the closed form of policy evaluation and then will converge rapidly, but if the dimension of the problem is higher (the number of action and states is very large) then the computation of the policy evaluation using the closed form will be computationally very expensive and then the PI will take much more time than the VI approach. The relative merits of the approches: -The pros of VI are: There isn't a lot of computation in each iteration (computationnally efficient). -The cons of VI are: The convergence is asymptotic and then the algorithm may need a lot of iterations to give a good approximation of the optimal value. -The pros of PI are: The PI approach converges to the exact optimal value unlike the VI approach in a finite number of iterations (generally not very large). -The cons of PI are: The PI approach is computationnally expensive since it needs to do at each iteration a complete policy evaluation which need a lot of calculation, specially when the matrix closed form is used.
In [2]:	Q4: denote with Qn(x,a) the value function estimated using Monte-Carlo, i.e., empirical average. Build such estimator and plot Jn - J π as a function of n from gridworld import GridWorld1 import gridrender as gui import numpy as np import time
	<pre>env = GridWorld1 ###################################</pre>
	<pre># print(env.state_actions[4]) -> [1,3] # print(env.action_names[env.state_actions[4]]) -> ['down' 'up'] # - env.gamma: discount factor ####################################</pre>
	# the action. Recall that in the terminal states only action 0 (right) is # defined. # In this case, you can use gui.renderpol to visualize the policy ####################################
	<pre>env.render = True state = 0 fps = 1 for i in range(5):</pre>
	<pre># here the v-function and q-function to be used for question 4 v_q4 = [0.87691855, 0.92820033, 0.98817903, 0.00000000, 0.67106071, -0.99447514, 0.00000000, - 0.82847001, -0.87691855,</pre>
	[0.00000000], [-0.62503460, 0.67106071], [-0.99447514, -0.70433689, 0.75620264], [0.00000000], [-0.82847001, 0.49505225], [-0.87691855, -0.79703229], [-0.93358351, -0.84424050, -0.93896668], [-0.89268904, -0.99447514]
	######################################
In [66]:	<pre>[1, 3] ['down' 'up'] def value_estimation(n,Tmax,gamma): Q=np.zeros((env.n_states,4)) N=np.zeros((env.n_states,4),dtype=int) for i in range(n): state=env.reset() if(0 in env.state_actions[state]): action=0 else:</pre>
	<pre>action=3 x_0,action_0=state,action N[x_0,action_0]+=1 t=0 S=0 term=False while(t<tmax action="0</pre" and="" env.state_actions[state]):="" if(0="" in="" nexts,="" reward,="" s+="gamma**(t-1)*reward" state="nexts" term="env.step(state,action)"></tmax></pre>
	else: action=3 t+=1 Q[x_0,action_0]+=S return np.divide(Q,N,out=np.zeros_like(Q), where=N!=0) The function value estimation takes in argument The number of iterations n, Tmax the maximal length of a trajectory and the discount_factor (gamma). It returns an array Q which is a 11x4 matrix and it has as rows the states and as columns the actions. The algo uses a deterministic action which is go right if it's possible, up otherwise.
In [3]:	<pre># here we build the function that calculate abs(Jn-J) at each iteration n v_q4 = [0.87691855, 0.92820033, 0.98817903, 0.000000000, 0.67106071, -0.99447514, 0.00000000, - 0.82847001, -0.87691855,</pre>
	<pre>J+=(1/11) *v[state] if(0 in env.state_actions[state]):</pre>
l	<pre>echantillon=[i for i in range(1000)] estimation=[plt_estimation(i,5,0.95,v_q4) for i in range(1000)]</pre>
	<pre>plt.plot(echantillon, estimation) plt.grid() plt.xlabel('iteration') plt.ylabel(' Jn-J ') plt.title('Figure2: Gap between the optimal value and the estimated one') plt.show()</pre> Figure2: Gap between the optimal value and the estimated one 0.30 0.25
	0.15 0.10 0.05 0.00 0 200 400 600 800 1000 iteration
	We could see that the gap decreases as the number of iteration increase with some fluctuations, the plot show the gap for a range of iteration from 0 to 1000. The time taken to plot the figure is not negligible duo to the expensive calculation in the function value_estimation as the number n becomes very large. But we could see that the gap is almost equal to 0 after the 10 000 th iteration (0.001). Q5: Describe the (parameters of the) exploration policy and learning rate chosen.
In [267]:	<pre>def q_learning(n, Tmax, gamma, epsilon): Q=np.zeros((11,4)) N=np.zeros((env.n_states,4), dtype=int) total_reward=0 for i in range(n): state=env.reset() random=np.random.rand() if (random<=epsilon):</pre>
	<pre>action = np.random.choice(env.state_actions[state]) else: z=Q[state,:][[env.state_actions[state]]] action=env.state_actions[state][np.argmax(z)] N[state,action]+=1 t=0 term=False while(t<tmax alpha="1/N[state,action]</pre" and="" nexts,="" reward,="" term="env.step(state,action)" total_reward+="reward"></tmax></pre>
	<pre>q=Q[nexts,:][[env.state_actions[nexts]]] Q[state,action]=(1-alpha)*Q[state,action]+alpha*(reward+gamma*np.max(q)) state=nexts if (random<=epsilon): action = np.random.choice(env.state_actions[state]) else: z=Q[state,:][[env.state_actions[state]]] action=env.state_actions[state][np.argmax(z)] N[state,action]+=1 t+=1 return [Q,total_reward]</pre>
	The function q_learning takes as arguments number of iteration n, the maximal length of the trajectory Tmax, gamma the discount factor and epsilon the exploration probability. It returns the 11x4 matrix of the Q values of each state*action and the total_reward accumulated over all the trajectories. The stepsizes is chosen to be harmonical (i.e α _N(xt,at) (state,action)=1/N[state,action] where N is the matrix given the number of times the (state,action) pairs are visited over all trajectories till the time t of the last trajectory. The exploration policy is chosen as parameter (argument) and will be chosen small enough to let all the state to be explorated.
In [407]:	<pre>illustrate the convergence of Q-Learning v_optimal =[0.877,0.928,0.988,0,0.824,0.928,0,0.778,0.824,0.877,0.828] iteration=[i for i in range(1000)] gap=[np.max(abs(v_optimal-np.max(q_learning(i,5,0.95,0.01)[0],axis=1))) for i in range(1000)] tot_reward=str(q_learning(1000,5,0.95,0.01)[1]) plt.plot(iteration,gap) plt.xlabel('iteration')</pre>
	plt.ylabel(' v-v* ') plt.grid() plt.suptitle('Figure3: Gap between the optimal value and the estimated one') plt.title('The total reward accumulated after 1000 iteration for ε=0.01 is: ' + tot_reward) plt.show() Figure3: Gap between the optimal value and the estimated one The total reward accumulated after 1000 iteration for ε=0.01 is: 617.0
	0.8 0.7 0.6 0.5 0.4 0 200 400 600 800 1000 iteration
In [361]:	<pre>v_optimal =[0.877,0.928,0.988,0,0.824,0.928,0,0.778,0.824,0.877,0.828] iteration=[i for i in range(1000)] gap=[np.max(abs(v_optimal-np.max(q_learning(i,5,0.95,0.1)[0],axis=1))) for i in range(1000)] tot_reward=str(q_learning(1000,5,0.95,0.1)[1]) plt.plot(iteration,gap) plt.xlabel('iteration') plt.ylabel(' v-v* ') plt.grid() plt.suptitle('Figure4: Gap between the optimal value and the estimated one') plt.title('The total reward accumulated after 1000 iteration for ε=0.1 is: ' + tot_reward)</pre>
	plt.title('The total reward accumulated after 1000 iteration for ε=0.1 is: ' + tot_reward) plt.show() Figure4: Gap between the optimal value and the estimated one The total reward accumulated after 1000 iteration for ε=0.1 is: 586.0
	v_optimal =[0.877,0.928,0.988,0,0.824,0.928,0,0.778,0.824,0.877,0.828] iteration=[i for i in range(1000)] gap=[np.max(abs(v_optimal-np.max(q_learning(i,5,0.95,0.4)[0],axis=1))) for i in range(1000)]
	<pre>gap=[np.max(abs(v_optimal-np.max(q_learning(i,5,0.95,0.4)[0],axis=1))) for i in range(1000)] plt.plot(iteration,gap) tot_reward=str(q_learning(1000,5,0.95,0.4)[1]) plt.xlabel('iteration') plt.ylabel(' v-v* ') plt.grid() plt.title('The total reward accumulated after 1000 iteration for ε=0.4 is: ' + tot_reward) plt.suptitle('Figure5: Gap between the optimal value and the estimated one') plt.show()</pre> Figure5: Gap between the optimal value and the estimated one The total reward accumulated after 1000 iteration for ε=0.4 is: 354.0
	0.6 ====================================
	v_optimal =[0.877,0.928,0.988,0,0.824,0.928,0,0.778,0.824,0.877,0.828] iteration=[i for i in range(1000)] gap=[np.max(abs(v_optimal-np.max(q_learning(i,5,0.95,1)[0],axis=1))) for i in range(1000)] plt.plot(iteration,gap) tot_reward=str(q_learning(1000,5,0.95,1)[1]) plt.title('The total reward accumulated after 1000 iteration for s=1 is: ' + tot_reward) plt.xlabel('iteration') plt.ylabel(' v-v* ')
	plt.ylabel(' v-v* ') plt.grid() plt.suptitle('Figure6: Gap between the optimal value and the estimated one') plt.show() Figure6: Gap between the optimal value and the estimated one The total reward accumulated after 1000 iteration for ε=1 is: -104.0
	abs (np.max (q_learning (5000, 5, 0.95, 0.1) [0], axis=1) -v_optimal)
Out[317]:	<pre>array([0.03035853, 0.01294132, 0.00263904, 0.</pre>
In [432]:	[0.87982636 0.92974869 0.98906333 0. 0.82569565 0.92360642 0. 0.77517404 0.72571866 0.86291688 0.69479574] It's maximum pointwise difference from the optimal value is (ϵ =0.1): 0.133204263896 value_es=np.max(q_learning(20000,5,0.95,0.9)[0],axis=1) print('The optimal value estimated one after 10000 iteration for ϵ =0.9 is:\n',value_es) print('\nIt\'s maximum pointwise difference from the optimal value is (ϵ =0.9):\n', np.max(abs(value_es-v_optimal))) The optimal value estimated one after 10000 iteration for ϵ =0.9 is:
	The optimal value estimated one after 10000 iteration for ϵ =0.9 is: [0.86274043 0.92378071 0.98694163 0. 0.78826138 0.92817264 0. 0.75349863 0.81891388 0.87945874 0.82865148] It's maximum pointwise difference from the optimal value is (ϵ =0.9): 0.0357386187317 By analyzing all those results, we could affirm that the value of ϵ represent the trade-off between the maximum reward and the convergence to the optimal value, by choosing ϵ very small, we don't allow for much exploration then the convergence to the optimal value will be very slowly w.r.t to the number of iterations but the accumulated reward will be significant (as shown in Figure3). If ϵ is very large then we allow for a lot of exploration, the convergence to the optimal value will come
In [78]:	<pre>def non_uniform_prior(n_state,r): z=np.random.random() if z<=r: state=np.random.randint(0,n_state//2+1) else: state=np.random.randint(n_state//2+1,n_state) return state</pre> The function non_uniform_prior generates a initial probability over states that is not uniform as like the function env.reset()
In [93]:	<pre>def value_estimation_bis(n,Tmax,gamma,r): Q=np.zeros((env.n_states,4)) N=np.zeros((env.n_states,4),dtype=int) for i in range(n): state=non_uniform_prior(env.n_states,r) if(0 in env.state_actions[state]): action=0</pre>
	<pre>else: action=3 x_0,action_0=state,action N[x_0,action_0]+=1 t=0 S=0 term=False while(t<tmax and="" nexts,="" reward,="" s+="gamma**(t-1)*reward</pre" state="nexts" term="env.step(state,action)"></tmax></pre>
	<pre>S+=gamma**(t-1)*reward if(0 in env.state_actions[state]): action=0 else: action=3 t+=1 Q[x_0,action_0]+=S return np.divide(Q,N,out=np.zeros_like(Q), where=N!=0)</pre> The function value_estimation_bis is the same as the function value_estimation where we change the prior distribution over
	The function value_estimation_bis is the same as the function value_estimation where we change the prior distribution over states from a uniform to a non_uniform one and we add an additionnal argument r: the parameters that control the non_uniform distribution over all states.

In [101]: Q_uniform_prior=value_estimation(10000,5,0.95)
Q_nonuniform_prior_1=value_estimation_bis(10000,5,0.95,0.3)
Q_nonuniform_prior_2=value_estimation_bis(10000,5,0.95,0.1)
Q_nonuniform_prior_3=value_estimation_bis(10000,5,0.95,0.9)

In [1]: import numpy as np
 import matplotlib.pyplot as plt
 import time

In [64]: P=np.zeros((3,2,3))

R=np.zeros((3,2)) P[0,0,0]=0.45

P[0,0,0]=0.45 P[2,0,0]=0.55 P[2,1,0]=1 P[2,0,1]=1 P[0,1,1]=0.5 P[1,1,1]=0.4 P[2,1,1]=0.1

P[2,0,1]=1 P[0,0,2]=0.6

Q1: Implement the discrete MDP model.