

Sessione 1

Nell'apprendimento per rinforzo (d'ora in poi verrà indicato con RL) si basa sul **processo decisionale di Markov** aka MDP attraverso un task di controllo dove un set di possibili stati e azioni ritornano un reward e una probabilità di passaggio ad uno stato all'altro. Il task nella pratica è un "compito" o una simulazione che per essere svolta (risolta) implica l'utilizzo dell'MDP.

NB: Il processo decisionale di Markov (MDP) asserisce che il passaggio allo stato successivo $T+1$, dipende esclusivamente dallo stato attuale T e non dagli stati precedenti. Quindi il processo **NON** ha memoria, in questo caso si dice che il processo è "**Markoviano**".

Tipologie di processi decisionali di Markov (MDP)

Esistono due tipologie di MDP, il processo a stati **finiti** e quello a stati **infiniti**.

Nel **MDP finiti** gli stati finiti hanno un numero finito di stati definiti dall'ambiente, es. l'uscita da un labirinto di 5 caselle x 5. In questo caso abbiamo 25 stati e 4 azioni (su, giù, dx e sx)

Nel **MDP a stati infiniti**, **invece**, l'ambiente appunto può restituire infiniti stati a fronte di infinite azioni, pensiamo per es. il sistema di guida automatica di una macchina dove l'azione es. girare il volante, è un valore continuo così come la scelta della velocità dell'auto.

Episodi

MDP definisce anche degli "episodi" in particolare:

Nel **MDP episodico** un episodio termina a determinate condizioni. Es. nel gioco degli scacchi quando la giocare da scacco matto.

Nel **MDP continuo**, il processo non ha fine, semplicemente continua ad esistere all'infinito in quanto non esiste uno stato fine.

Traiettoria ed episodio

La **traiettoria** è il movimento che l'agente compie per muoversi da uno stato all'altro. La **traiettoria** è definita dal simbolo greco τ "tau". Un esempio può essere:

$\tau = S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3$ che indica la partenza dallo stato zero, dove viene effettuata l'azione A_0 che porta la reward R_1 e il posizionamento nello stato S_1 , alla quale segue l'azione A_1 e così via. L'ultimo stato della traiettoria sarà (nel caso specifico) S_3 .

La **traiettoria** può essere **finita** o **infinita**, se è finita si chiama **episodio** è semplicemente una traiettoria che inizia in uno stato e finisce nello stato finale oltre quale non si torna indietro. es:

$\tau = S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3, \dots, \mathbf{R_T, S_T}$.

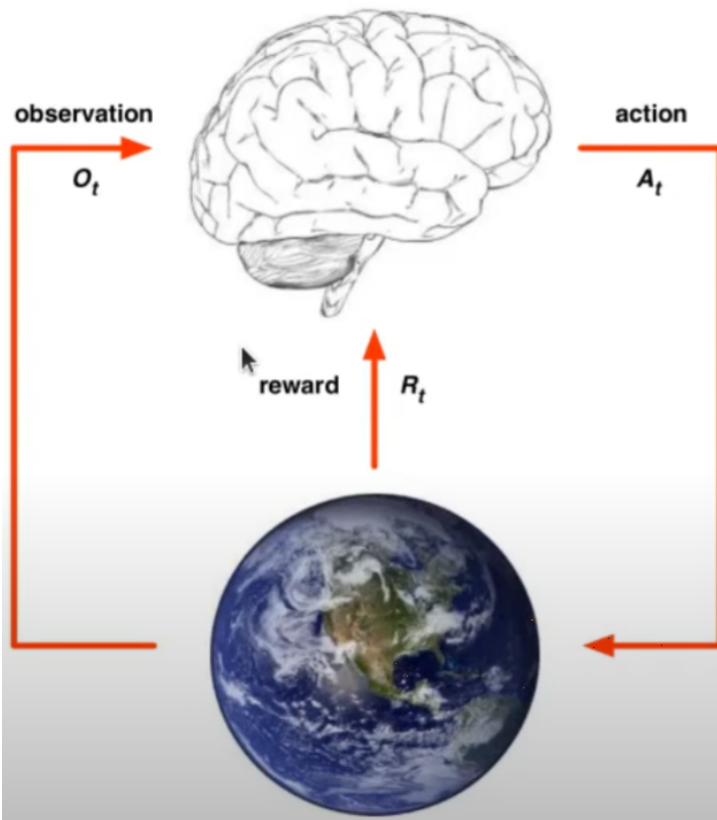
We are not alone! Second actor: the *environment*

Agent, step t

- Receives observation O_t .
- Receives scalar reward R_t .
- Computes his own *state* S_t^a .
- Executes action A_t .

Environment, step t

- Receives action A_t .
- Computes his own *state* S_{t+1}^e .
- Emits observation O_{t+1} .
- Emits scalar reward R_{t+1} .



di cui deriva il concetto di storia, ovvero la somatoria delle osservazioni, ricompense e azioni fino all'azione finale.

Notation

- **History**: the sequence of observations, actions, rewards up to time step t :

$$H_t := O_1, R_1, A_1, \dots, A_{t-1}, O_t, R_t.$$

- The agent selects actions, and the environment answers with *observations* and *rewards*.
- **State**: the information used (by the agent and the environment) to determine what happens next.
- State is naturally a sequence S_t .
- Agent state is a function of history: $S_t := f(H_t)$.
- **Environment state** S_t^e is different from **agent state** S_t^a .

Ricompensa vs Ritorno

La **ricompensa** (*reward*) viene restituita a fronte di un'azione, quindi per risolvere un task, dobbiamo massimizzare le ricompense ottenute. La ricompensa è quindi un risultato **immediato**. (R_t)

Invece il **ritorno** è la somma di tutte le ricompense ad un determinato momento nel tempo es: **G(t)** = $R(t+1) + R(t+2) + \dots + R(T)$ finchè il task è stato completato.

Predizione, miglioramento e Controllo

La predizione significa calcolare il valore di una certa policy fissata, il miglioramento, come dice la parola serve per migliorare (anche solo di poco, non la migliorare in assoluto) la policy attuale, mentre il controllo serve per trovare la migliore di tutte.

Da qui il ciclo utile per trovare la policy migliore: $\pi(s) \rightarrow V\pi(s) \rightarrow \pi' \geq \pi$ che inserire in un loop fino a quando converge. (teorema banach caccioppoli)

Fattore di sconto γ (gamma)

Il fattore di sconto è un incentivo per completare l'episodio nel miglior modo (più efficiente) possibile. Per ottenere questo la ricompensa dovrà essere moltiplicata per il fattore di sconto che diminuirà nel tempo all'aumentare delle azioni intraprese, rendendo le ricompense sempre più basse e quindi disincentivando le traiettorie lunghe.

Il fattore è un valore compreso tra zero e uno e viene elevato ad un esponente corrispondente dall'iesima azione fino alla fine dell'episodio.

Se γ (gamma) vale zero l'agente cercherà di prendere una ricompensa immediata, il che denota una strategia miope che non ottimizza l'apprendimento. Al contrario invece, un fattore γ gamma pari a uno, rende l'agente più "paziente" e quindi non prono ad ottimizzare gli step dell'episodio. In genere il fattore gamma viene settato a 0,99 che forza l'agente ad avere una ricompensa immediata ma allo stesso tempo lo forza ad avere una visione "a lungo termine".

In conclusione il fattore gamma indica all'agente quanto può valutare in maniera ottimali le azioni future.

NB: se il task è *episodico* allora è possibile utilizzare un fattore di sconto pari a 1, diversamente se il task è *continuo* (infinito senza stati assorbenti terminali) allora il tasso di sconto è meglio che sia <1 .

Policy

La policy dell'agente è una funzione che prende in input uno stato e ritorna l'azione che va presa in quello stato. E' rappresentata dalla lettera greca π

$$\pi : S \mapsto A$$

La probabilità di eseguire un'azione (a) nello stato (s) si può rappresentare come: $\pi(a|s)$

L'azione che la policy sceglie nello stato (s) viene descritta dalla formula: $\pi(s)$

Dipende quindi dal contesto, in alcuni casi si utilizza il primo $\pi(a|S)$ in altri il secondo $\pi(s)$.

La policy può essere di due tipi: **stocastica** o **deterministica**.

Si dice che la policy è **deterministica** quando viene scelta **sempre** la stessa azione in un determinato stato quindi stiamo parlando del caso $\pi(S)$.

Si dice invece **stocastica** quando l'azione viene scelta sulla base delle probabilità es. : $\pi(S) = [0.3, 0.2, 0.5]$ ovvero la probabilità di effettuare una azione nello stato S, è del 30% nel primo caso, 20% nel secondo e 50% nel terzo. Quindi siamo in presenza del caso $\pi(a|S)$

Quindi bisogna trovare la policy ottimale rappresentata come π^* (greco - asterisco) che sceglie le azioni che massimizzano la somma dei fattori di sconto per le ricompense alla lunga.

La policy π è una distribuzione di probabilità che decidiamo noi dato lo stato con la quale scegliamo le azioni e prendendo quindi una decisione, si differenzia dalla probabilità che NON decidiamo noi che si rappresenta con \mathbf{P} che invece rappresenta il modello la cui probabilità non possiamo modificare. (es. $P(s',r|s,a)$)

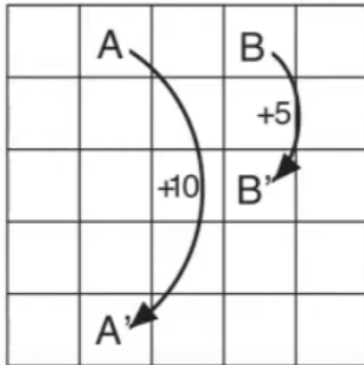
Il valore della policy V in genere indica la ricompensa totale media che si ottiene applicando la policy

Controllo Ottimale

Nel RL è fondamentale determinare la Policy Ottimale π^* indispensabile per la gestione dell'ambiente. La griglia di centro (visibile in figura) è la rappresentazione del valore di ciascuno stato. (dove per stato si intende ogni casella della griglia) Per valore ottimale V^* si intende il valore della policy ottimale, ovvero quello che si può ottenere facendo le azioni migliori possibili. Per migliore azione si intende determinare lo scopo, ovvero massimizzare le somme delle ricompense (dette *ritorno*) ottenibili con le azioni future.

La policy ottimale quindi, si ottiene valutando di volta in volta il valore ottimale. Le policy ottimali sono tante, anche su un unico stato, vedi per es. che nella casella in fondo a sx il cui valore è 14,4 ha due policy ottimali in quanto i valori ottimali in questo specifico caso sono due.

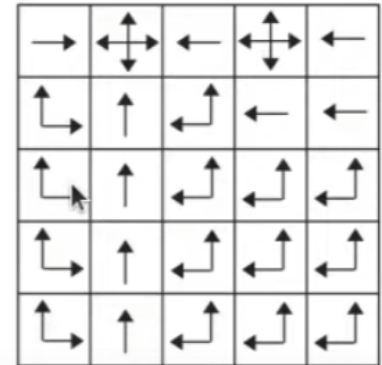
Gridworld example: *optimal control*



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

b) v_*



c) π_*

Exercise

- Compute the *optimal value* function over all possible policies.
- Given the optimal value v_* as above, find the optimal policy.
- Is the optimal policy unique?

Pianificazione e Apprendimento

L'apprendimento nel RL si basa sulla pianificazione.

La *pianificazione* implica la conoscenza del modello associato all'ambiente, es. il lancio di un dado che implica che il *valore medio detto anche ritorno medio* dell'azione è 3,5 ovvero $1*1/6+2*1/6+3*1/6+4*1/6+5*1/6+6*1/6$.

NB: il *modello* è l'ambiente e normalmente NON lo conosciamo.

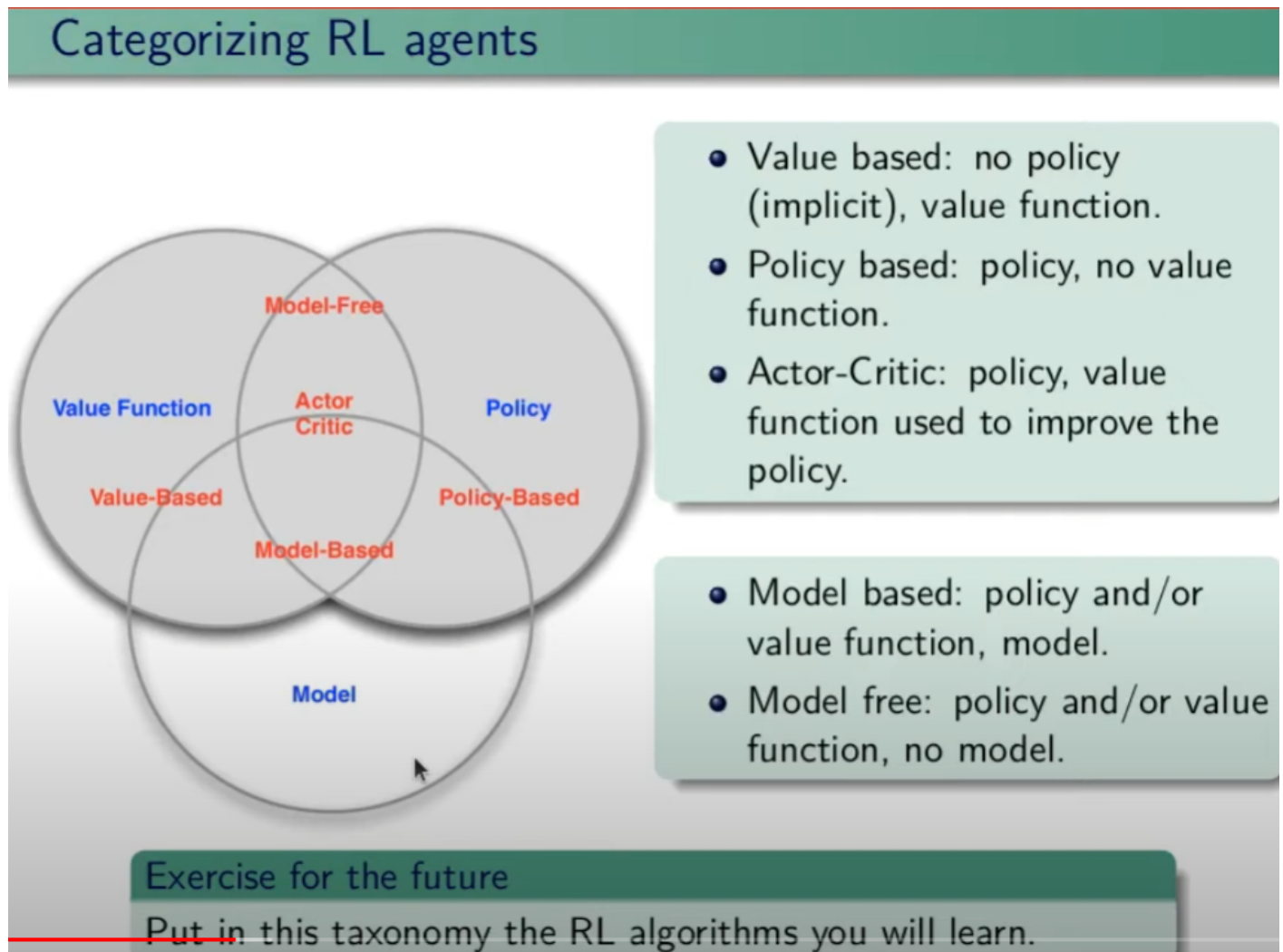
L'apprendimento, è la fase successiva alla pianificazione e implica *l'iterazione con l'ambiente* e prevede il calcolo della *media empirica* ovvero la media dei valori ottenuti dall'iterazione con l'ambiente.

Quindi con la pianificazione e l'apprendimento, l'agente migliora la policy.

Entrambi guardano avanti nel futuro calcolando i valori cercando il miglioramento della policy.

Tipologie di algoritmi applicabili al RL

Di seguito viene rappresentata la tassonomia (categorizzazione) delle varie tipologie di algoritmi applicabili nell'ambito del RL.



Esplorazione vs Sfruttamento

E' l'eterno dilemma, ovvero provo sempre nuove soluzioni o sfrutto sempre quelle che già conosco. Entrambi vanno utilizzati per "allenare" la rete neurale. (vedremo più avanti)

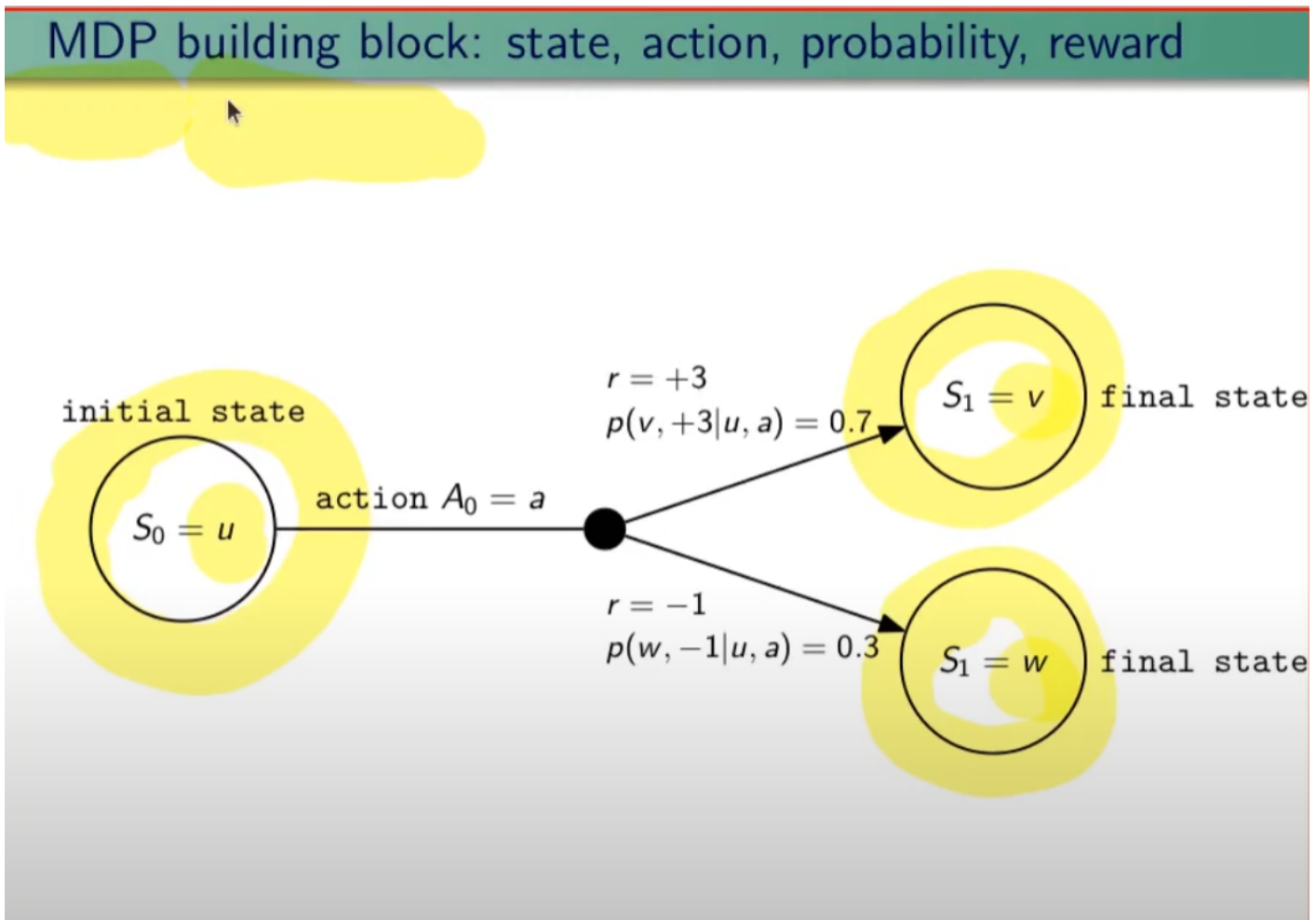
Processi decisionali di Markov (MDP)

I processi decisionali di Markov si basano sulla gestione degli stati. Gli stati vengono rappresentati come dei cerchi dove all'interno è presente una lettera che lo rappresenta. Il pallino nero invece è l'azione.

Nello schema sotto rappresentato vediamo che, partendo dallo stato "u", la risposta dell'ambiente a fronte dell'azione "a" di tipo probabilistico (aleatoria) e quindi non deterministico, è di 0,7 con un reward di +3 che ci porta di v, oppure 0,3 con reward -1 che ci porta in "w". Se invece a fronte di una azione (pallino) ci fosse stato una unica determinazione della risposta dello stato T+1 anzichè

due o piu, allora la risposta sarebbe stata derministica.

Es. leggendo la risposta dell'ambiente a fronte dell'azione a (nel caso dello 0,7 - 70%) è: la probabilità di arrivare in "v" con una ricompensa di +3, dato che sono partito dallo stato iniziale "u" e ho fatto l'azione "a".



Sempre rimandando nel diagramma sopra riportato, qual'è la probabilità di andare in "w" con una reward di 42? (vedi figura sotto)

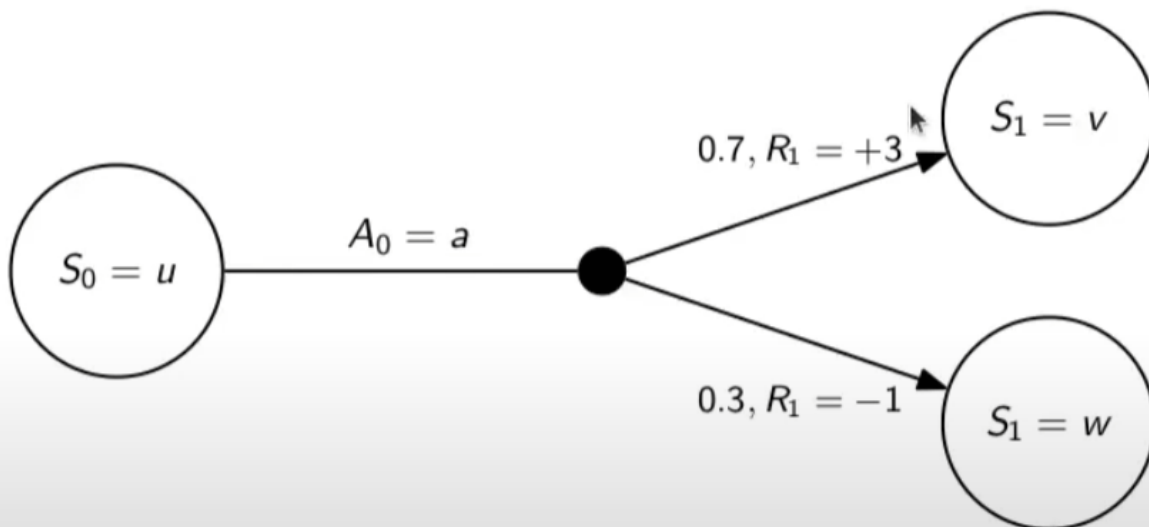
$$p(w, 42|u, a) = ?$$

La risposta è zero in quando deve verificarsi la combinazione stato + reward che in questo caso non esistendo è quindi pari zero.

NOTA: lo **stato terminale** si indica con un quadrato.

Ora semplifichiamo un attimo lo schema per renderlo più leggibile come sotto riportato:

MDP building block: state, action, probability, reward



Exercise

Discuss in details what S_0 , A_0 , S_1 and R_1 are.

Nello specifico la differenza tra S_0 e S_1 indica la situazione dello stato "S" al tempo T. Ovvero S_0 è lo stato iniziale, mentre S_1 è sempre lo stato S ma al tempo T+1 e può valere "w" o "v" a seconda della probabilità. Se ci fosse uno stato successivo S_2 dovremmo considerare le probabilità di transizione da S_0 a S_1 e da S_1 a S_2 .

La rappresentazione è quindi un **grafo orientato** in quanto unidirezionale. i tondi sono le azioni e gli archi i risultati probabilistici delle azioni.

Ho quindi un insieme di stati S e di azioni A e di rewards R. Dopo ogni azione ho quindi una distribuzione di probabilità su $S \times R$. Per ogni coppia di stato azione ho una distribuzione di probabilità. Conoscere il modello significa quindi conoscere la distribuzione di probabilità. Introduciamo anche il fattore di sconto gamma che vale sempre di meno all'allungarsi della distanza di tempo in modo che l'agente sia incentivato a trovare la soluzione migliore più

velocemente.

La funzione p rappresenta quindi il modello, è quindi la probabilità di transizione dallo stato al tempo $T-1$ allo stato al tempo T . (a fronte di una certa azione e relativa ricompensa)

Generalizzando quanto sinora detto:

Markov Decision Process: meaning of the model

From distribution model to random variables S_t and R_t

The probability distribution p of the MDP gives the *next* state and reward:

$$\Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a) := p(s', r | s, a).$$

Markov decision process data

- A set of *states* \mathcal{S} , a set of *actions* \mathcal{A} and a set of *rewards* \mathcal{R} .
- For each state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$, a probability distribution $p(\cdot, \cdot | s, a)$ over $\mathcal{S} \times \mathcal{R}$.
- A discount factor $\gamma \in [0, 1]$.

“ dove il simbolo \square (pipe) indica che un elemento "a" è condizionato **da** "b", es. "a | b". Occhio che per es. $P(a,b)$ è la probabilità dato in input "a" e "b" mentre invece $P(a|b)$ significa la probabilità dato "a" e condizionato a "b" MA "a" può assumere più valori in B.

dove il simbolo \sim (tilde) indica che uno stato S è preso da una certa distribuzione di probabilità P , es. $S \sim P(\dots)$

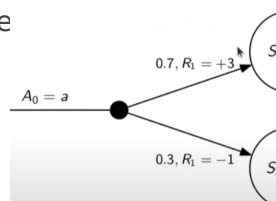
Esercizio:

Exercises

- Explain what S_t , A_t and R_t are.
- Given p , give a formula for $\Pr(S_t = s' | S_{t-1} = s, A_{t-1} = a)$.
- Given p , give a formula for $\mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a]$.
- Given p , give a formula for $\mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s']$.

nota: ricordare che il simbolo \mathbb{E} rappresenta il ritorno atteso ovvero il valor medio dei ritorni.

In pratica \mathbb{E} rappresenta il valore medio (**media dei valori pesati**) delle ricompense che posso ottenere in un determinato stato S . Questo perchè in uno stato potrei avere una distribuzione di probabilità che per es. a fronte di una azione



probabilità e un ricompensa ciascuna. es:

1)

$$\Pr(S_t = s' | S_{t-1} = s, A_{t-1} = a).$$

La prima riga si legge così: la probabilità che un stato S al tempo T sia s' condizionato al fatto che lo stato precedente fosse s e che io abbia fatto precedentemente una decisione azione a . Quindi si vuole sapere la probabilità di essere nello stato s' nel caso in cui nello stato precedente s sia stata eseguita l'azione a .

Per ricavare questa probabilità bisogna partire la formula generica MDP che tiene conto delle rewards, ovvero:

$$\Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a) := p(s', r | s, a).$$

quindi per rispondere al primo quesito dobbiamo sommare tutte le rewards in modo che rimanga la sola probabilità che ci porta allo stato s' , in quanto non vogliamo la probabilità che esca s' e r , ma vogliamo solo la probabilità che esca s' , ovvero:

$$P(s', r | s, a) = P_2(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

↳

$$P_2(S_t = s' | S_{t-1} = s, A_{t-1} = a)$$
$$= \sum_{r \in \mathcal{R}} P(s', r | s, a)$$

perchè sommando **tutte** le rewards ho la certezza di finire nello stato s' . Sommare tutti gli "r" si dice anche saturare tutti gli indici r.

2)

$$\mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a].$$

In questo caso invece si desidera la ricompensa media (variabile aleatoria) al tempo "t" condizionata dal fatto che al tempo t-1 partissi dallo stato "s" e facessi l'azione "a".

Quindi il valore medio della ricompensa sarà la sommatoria di tutte le ricompense ciascuna di esse moltiplicate per la propria probabilità di uscita, ovvero:

$$= \sum_{r \in R} r \sum_{s' \in S} \cdot p(s', r | s, a)$$

dato che in questo caso vogliamo la reward media al tempo t , come nel caso 1) dobbiamo "saturare" un fattore, in questo caso sono gli stati. Quindi per tutti gli stati s' moltiplichiamo la probabilità ritornata dall'azione a nello stato " s " a $t-1$ per tutte le reward. La prima sommatoria estrare tutte le reward che vanno a moltiplicare le corrispettive probabilità. Moltiplicando reward per la probabilità si ottiene la reward media. (vedi esempio del valore medio del lancio del dado, dato dalla sommatoria del valore di ciascuna faccia del dado per $1/6$ per un totale di $3,5$)

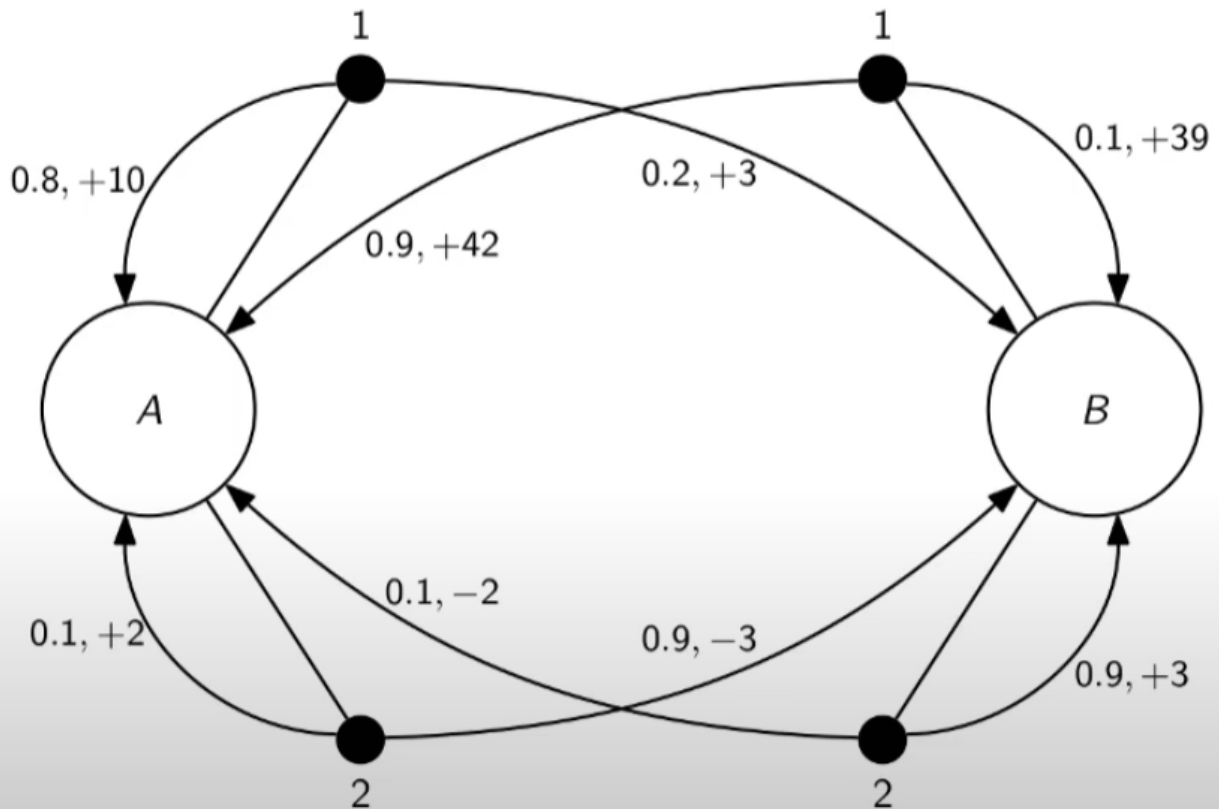
3) per ora non viene spiegato.

Dinamica del MDP: rappresentazione tabellare

Una delle tecniche più basiche di rappresentazione delle transizioni da uno stato all'altro fa uso di matrici di valori che indicano la probabilità di transizione da uno stato all'altro. Se per esempio un sistema ha 4 stati, la tabella sarà composta da una matrice di 4×4 . Ovviamente questo metodo funziona con sistemi i cui stati sono molto limitati, non per funziona con sistemi complessi come per es. gli scacchi o la dama.

Facciamo un esempio, calcoliamo su questo MDP delle quantità:

The example MDP



Iniziamo con il calcolo di una matrice di transizione detta matrice stocastica, in questo caso avendo due stati la matrice è 2x2, nelle celle della matrice andremo ad inserire le probabilità di transizione da uno stato all'altro. Nell'esempio sotto riportato andremo a calcolare la matrice associata all'azione 1. La peculiarità di questa matrice è che la somma di ciascuna riga da sempre 1 - 100%.

P = 1 (azione 1)	A	B
A	0,8	0,2
B	0,9	0,1

Calcoliamo ora il vettore ricompensa associata all'azione 1.

R = 1 (azione 1)	$10 \cdot 0,8 + 3 \cdot 0,2 = \mathbf{8,6}$	$0,9 \cdot 42 + 0,1 \cdot 39 = \mathbf{41,7}$
------------------	---------------------------------------------	-----------------------------------------------

ora creiamo che tabella che rappresenta l'intero modello:

s (stato partenza)	a (azione)	s' (stato di arrivo)	r (ricompensa)	p(s',r s,a) probabilità
A	1	B	3	0,2
...				

--	--	--	--	--

La tabella conterrà un totale di righe dato dalle azioni x numero di probabilità di accedere allo stato, questo caso 8.

La proprietà di "Markovianità"

La proprietà di markov indica che il valore dello stato St dipende esclusivamente dallo $St-1$ ma **non** dipende dagli stati precedenti a $St-1$. (es. $St-2$, $St-3$ etc etc)

Episodi MDP

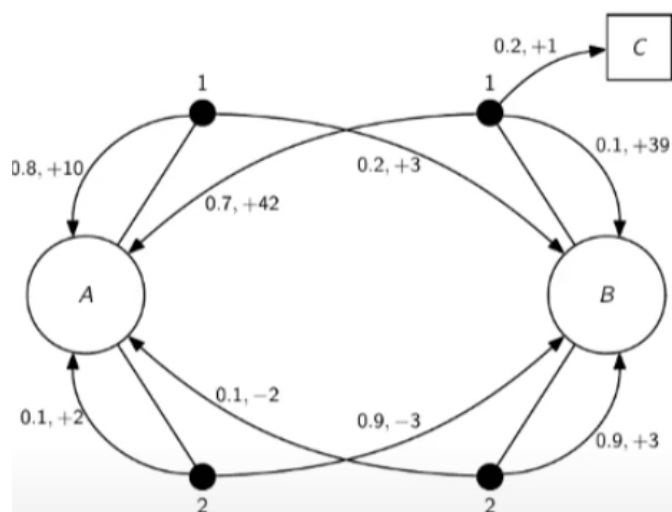
Se ho uno stato terminale tra gli stati possibili allora il mio task si chiama *episodico*. Attenzione che però dipende dalla policy, ovvero se ho delle azioni che mi consentono di evitare lo stato terminali posso entrare in loop. E' compito dell'algorithm evitare questi loop. Nel caso in cui non esista lo stato terminale si tratta di un *MDP continuing*.

Esistono anche task a *orizzonte temporale definito* se esistono dei limiti temporali.

Un episodio quindi è una sequenza di stati, azioni, rewards che terminano con lo stato terminale.

Simuliamo un episodio:

Episodic MDP



- If there are special *final states* reachable from every state, the MDP is *episodic*.
- Otherwise, the MDP is *continuing*.
- *Episode*: any sample $S_0, A_0, R_1, S_1, \dots$ terminating in the final state.

Exercise

- Write an episode, and compute its probability of happening. Hint: tricky question.

EPISODIO: A, 1,10, A, 2, -3, B, 2, 0,9, B, 1, 1, C

PROBABILITA': dipende da:

- probabilità di scegliere (1|A)
- probabilità di scegliere (2|A)
- probabilità di scegliere (1|B)
- probabilità di scegliere (2|B)

Per calcolare un episodio bisogna capire quindi come viene definita la policy.

Un esempio di policy π che per esempio può essere scegliere sempre l'azione 1. Quindi la policy "pi" è fondamentalmente una strategia che può (e spesso deve) variare per massimizzare il ritorno.

Ritorno

Il ritorno ($G = \text{gain}$) è la somma delle ricompense di un episodio ed è una variabile aleatoria.

The *return*: towards the goal

Definition

- The *total return* of an episode ending at time T is the value of the random variable

$$G_t := R_{t+1} + R_{t+2} + \dots + R_T$$

for that episode.

- If the MDP is continuing, we need a *discount factor*:

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

Why?

- In episodic tasks we use $\gamma = 1$.

NB: L'ipotesi è che i ritorni R abbiano un limite superiore finito che consente al fattore di sconto far convergere il ritorno totale.

The return: towards the goal

Why the discount

- The discount factor measures how much do we care about rewards far in the future.
- A reward r after $k + 1$ time-steps is worth "only" $\gamma^k r$: we say *myopic evaluation* if $\gamma \sim 0$, *far-sighted evaluation* if $\gamma \sim 1$.
- Convenience: avoids infinite returns in cyclic MDP.
- We shouldn't trust our model too much: uncertainty about the future may not be fully represented.
- If the reward is financial, immediate rewards may earn more interest than delayed rewards.
- Animal and human behaviour shows preference for immediate reward.

$$P_{ss'}^a$$

Ricordo la notazione che si ripeterà spesso nel corso, $P_{ss'}^a$ ovvero la probabilità che facendo un'azione "a" nello stato "s" l'ambiente mi porti nello stato s'. Il mio comportamento è rappresentato da una distribuzione di probabilità negli stati S. La strategia di scelta delle azioni - **che sia chiama policy** - è una distribuzione di probabilità negli stati.

Le policy ottimali possono essere deterministiche o non deterministiche. Un esempio di policy deterministica è il lancio de dado in quanto ho una sola possibile azione, appunto il lancio del dato, invece un esempio di non deterministica è per es. sasso-carta-forbice.

Funzione stato-valore

La funzione stato valore per un MDP è il numero che ci indica quale azione intraprendere. Si rappresenta come:

Definition: state-value function

The *state-value function* $v_{\pi}(s)$ for a MDP is the return we can expect to accumulate starting from state s , following the policy π :

$$v_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s]$$

dove V è il valore indicizzato dalla lettera π (policy) indicizzato dallo stato S .

Mentre "**E**" è il **valore atteso/valore medio**, è il valore dello stato, è il numero teorico che se lo conoscessimo aprioristicamente ci "renderebbe la vita più semplice" in quanto significherebbe che conosceremmo l'ambiente, cosa ovviamente non possibile. (spesso il valore E è la media ponderata dei valori attesi)

"**E**" il è il valore che voglio massimizzare, viene quindi indicato con la lettera "**V**" che sta per valore indicizzato a π -greco ovvero la "policy" cioè quello che facciamo noi espresso in probabilità di effettuare un'azione, dello stato " S " (notazione funzionale) uguale al valore medio atteso di G_t -> ritorno al tempo " t " condizionato al stato " s " che è l'argomento della funzione $V_{\pi}(s)$

Q-learning

E ora veniamo ad un concetto molto importante, il concetto del **Q-learning**, ovvero action value function che sta per qualità dell'azione.

$Q_{\pi}(s,a)$ = rappresenta la **qualità di una azione** -> partendo dalla prima azione (**fissata**) e dallo stato (**fissato**) eseguo tutte le azioni che la policy mi dice di eseguire. Tenendo ovviamente in conto che il modello (P) mi dice quale sarà il nuovo stato e la ricompensa e la policy π deciderà quale sarà la nuova azione da eseguire.

$V_{\pi}(s)$ = rappresenta il **valore dello stato** " s " (stato di partenza) che si ottiene seguendo le azioni dettate dalla policy π tenendo conto che la scelta delle azioni è di tipo probabilistico quindi in teoria sarebbe πP . (" P " si omette) Questo guadagno è una variabile aleatoria. L'intenzione è quello di avere il valore medio (la quantità) più grande possibile

$V^*(s)$ = è il valore medio massimo ottimale che si può ottenere su tutti i $V_{\pi}(s)$ applicando tutte le policy possibili. -> $\max_{\pi} V_{\pi}(s)$

$Q^*(s,a)$ = è la funzione valore ottimale -> $\max_{\pi} q_{\pi}(s,a)$

In fine definiamo la **policy ottimale** π^* in grado di massimizzare il valore (sia l'azione valore che lo stato valore)

The **action-value function** $q_\pi(s, a)$ is the expected return from a state s , choosing action a , and then **following the policy** π :

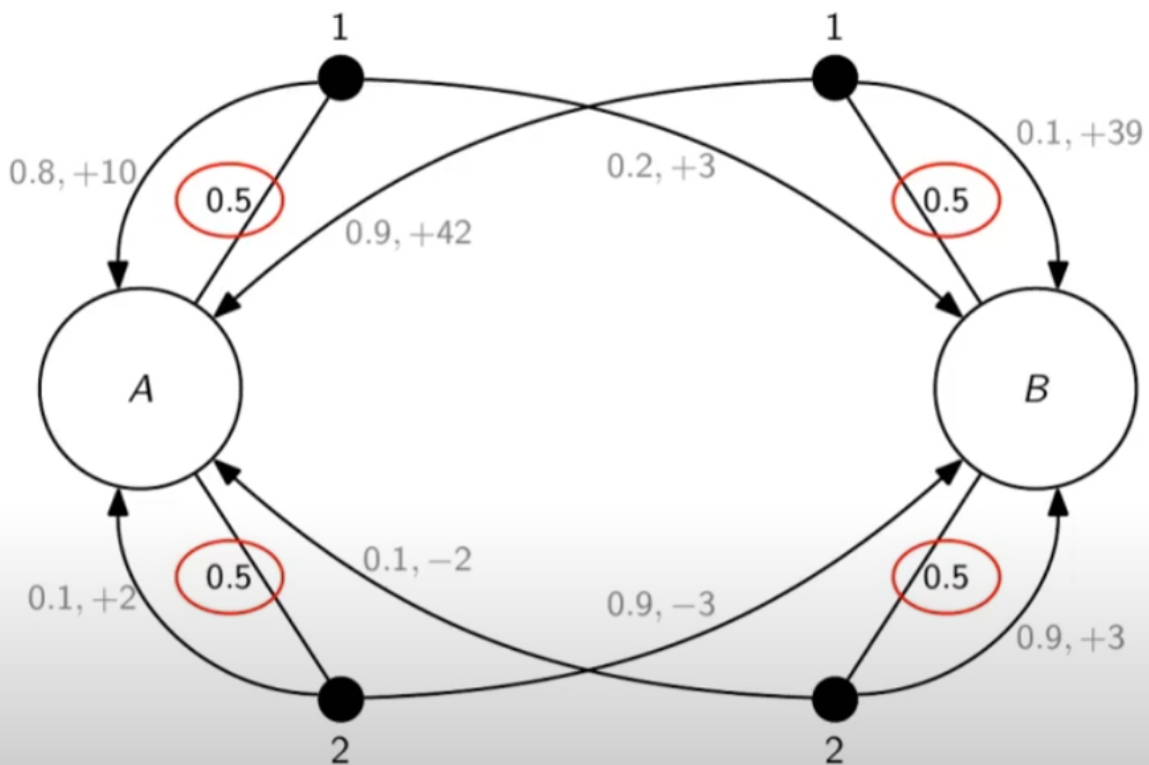
$$q_\pi(s, a) := \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

Similarly, the **state-value function** is:

$$v_\pi(s) := \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_{a \sim \pi}[q_\pi(s, a)]$$

Ora facciamo un esempio:

Example



Exercise

Compute $q_\pi(B, 1)$, $q_\pi(B, 2)$ and $v_\pi(B)$ for the uniform policy π .

Esercizio: calcolare il q-values applicando la policy π nello stato B e compiendo l'azione 1 in un caso e l'azione 2 nell'altro; e calcolare il valore atteso nello stato B applicando la policy uniforme π .

Soluzione $q_{\pi}(B,1)$

$$q_{\pi}(B,1) = \gamma * 0,1 * (39 + V_{\pi}(B)) + \gamma * 0,9 * (42 + V_{\pi}(A))$$

dove:

$$V_{\pi}(B) = 0,5 * q_{\pi}(B,1) + 0,5 * q_{\pi}(B,2)$$

$$V_{\pi}(A) = 0,5 * q_{\pi}(A,1) + 0,5 * q_{\pi}(A,2)$$

spiegone:

Per determinare la funzione q (valore azione) relativa alla policy π che nel caso di partenza da B, è necessario applicare la "policy di partenza" ovvero al 50% di probabilità di scegliere un'azione. Quindi la $q_{\pi}(B,1)$ è calcolata come la probabilità di scegliere l'azione 1 (che in questo caso è del 100% in quanto viene impostato come richiesta iniziale) che moltiplica la probabilità del 10% di ottenere 39 sommato al valore- π dello stato di arrivo ovvero "B", più la probabilità del 90% di andare nello stato A ottenuto 42 sommato al valore- π di A.

Il valore- π di A è a sua volta la probabilità di scegliere l'azione 1 sommata alla probabilità di scegliere l'azione 2 ovvero: $V_{\pi}(A) = 0,5 * q_{\pi}(A,1) + 0,5 * q_{\pi}(A,2)$. (dove la probabilità per questa policy di scegliere l'azione è del 50%)

Analogamente il valore- π di B è probabilità di scegliere l'azione 1 sommata alla probabilità di scegliere l'azione 2 ovvero: $V_{\pi}(B) = 0,5 * q_{\pi}(B,1) + 0,5 * q_{\pi}(B,2)$. (dove la probabilità per questa policy di scegliere l'azione è del 50%)

Da notare che questa non è la policy ottimale in quanto è la policy di partenza che ci da il 50%, l'apprendimento sta **nell'iterare** il processo a fine di variare la % della policy per farla convergere a quella ottimale.

Bisogna quindi calcolare 4 equazioni con 4 incognite che sono: $q_{\pi}(A,1)$ $q_{\pi}(A,2)$ $q_{\pi}(B,1)$ e $q_{\pi}(B,2)$

Per farla convergere però, bisogna utilizzare un fattore di sconto gamma, altrimenti tende a infinito... (di qui l'introduzione di gamma)

Questo ciclo (iterazione) viene anche detto **controllo**.

In conclusione questo sistema è una versione dell'equazione di Bellman.

NB: per ricompensa si intende la somma dei valori che il modello ci da quando da uno stato passa ad un altro. Quindi la ricompensa è composta da uno a più valori con relativa percentuale di

ottenimento. (es. mi da valore 20 al 10%, valore 50 al 60% e così via) Una volta ottenuti tutti i valori dello stato si calcola il valore medio che appunto rappresenta la ricompensa.

Equazioni di Bellman

Ed eccoci ad uno dei capisaldi dell'apprendimento per rinforzo, le equazioni di Bellman che si basano sui concetti sinora appresi, ovvero ritorno (somma delle ricompense), ricompensa, valore dello stato, funzione valore-azione.

Exercise

Find a formula for the value function $v_{\pi}(s)$ in a certain state s as a function of $v_{\pi}(s')$, for successors s' of s . Hint: first prove that

$$G_t = R_{t+1} + \gamma G_{t+1}$$

Il ritorno G_t può essere formulato (descritto) in maniera **ricorsiva** in quanto è la somma delle ricompense.

Ovvero il ritorno al tempo t è la somma della ricompensa al tempo $t+1$ sommato al ritorno del tempo $t+1$, dove il ritorno al tempo $t+1$ è rispetto a $t+1$ (quindi es. $t=2$ scontato di gamma ovviamente)

dimostrazione:

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) \\ &\quad \underbrace{\hspace{10em}}_{G_{t+1}} \end{aligned}$$

quindi

$$G_t = R_{t+1} + \gamma G_{t+1}$$

cvd

da qui ne deriva l'equazione di Bellman:

Funzione valore dello stato

La **prima** equazione ricorsiva di Bellman

Theorem: state-value Bellman equation

The state-value function satisfy the following recursive formula:

$$v_{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{\pi}(s') \right)$$

innanzitutto bisogna osservare che è un'equazione **ricorsiva**, in quanto all'interno della funzione è riproposta la funzione stessa.

Come si legge l'equazione di Bellman:

Partiamo dal presupposto che a questo livello stiamo **pianificando**, ovvero cerchiamo di capire cosa potrebbe accadere se facessimo l'azione a_1 piuttosto che a_2 . In futuro vedremo come apprendere...

Quindi:

Partiamo dallo stato "s", abbiamo una policy π che indica la probabilità di scegliere un'azione tra le azioni

$$\sum_{a \in \mathcal{A}} \pi(a|s)$$

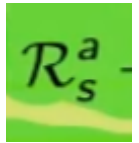
disponibili, con una certa probabilità, ovvero ->



$$\left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{\pi}(s') \right)$$

e quindi vediamo cosa succede, ovvero ->

che significa:

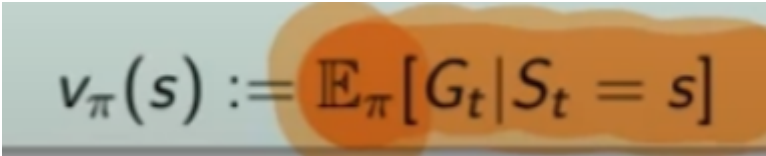
la ricompensa media (è un singolo numero) ottenuta dallo stato "s" facendo l'azione "a"


$$R_s^a$$

sommata al valore dello stato s',  la cui probabilità di arrivarci è data da  che rappresenta quindi la probabilità di passare dallo stato s a s', ed è una matrice di valori detta anche *matrice di transizione*.

dimostrazione:

partiamo dall'assunto:


$$v_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s]$$

che ci dice che il valore-pi allo stato S è pari al valore atteso dato dal ritorno al tempo t subordinato allo stato "s".

$v_{\pi}(s) =$

$$\begin{aligned}
 v_{\pi}(s) &= \mathbb{E}_{\pi} [G_t | S_t = s] = \\
 &\mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= \mathbb{E}_{\pi} [R_{t+1} | S_t = s] + \\
 &\quad \gamma \mathbb{E}_{\pi} [G_{t+1} | S_t = s] \\
 &= \sum_a \pi(a|s) \cdot \mathbb{E}_{\pi} [R_{t+1} | S_t = s, A_t = a] \\
 &\quad + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_t = s, A_t = a]
 \end{aligned}$$

+

$$+ \gamma \mathbb{E}_{\pi} [G_{t+1} | S_t = s, A_t = a]$$

ovvero:

$$\mathbb{E}_\pi [R_{t+1} | S_t = s, A_t = a] = \mathcal{R}_s^a$$

$$\mathbb{E}_\pi [G_{t+1} | S_t = s, A_t = a] = \sum_{s'} \mathcal{P}_{ss'}^a V_\pi(s')$$

Funzione valore stato-azione

La seconda equazione ricorsiva di Bellman

Theorem: Bellman equation

The action-value function satisfy the following recursive formula:

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a'|s') q_\pi(s', a')$$

Exercise

Prove the action-value Bellman equation.

Valori migliori possibili delle stati e delle azioni (policy ottimali V^* e Q^*)

Per migliorare la policy devo cercare il "meglio possibile" sia degli stati-valore che degli stati-valore-azione nello stato s per tutte le policy possibili, e quindi sarà il valore più alto (max) ottenibile. Questi valori saranno indicati come $V^*(s)$ e $Q^*(s,a)$

Definition

The *optimal state-value function* v_* is the maximum state-value over all policies:

$$v_*(s) := \max_{\pi} v_{\pi}(s)$$

Definition

The *optimal action-value function* q_* is the maximum action-value over all policies:

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

Però bisogna dire che le policy potrebbero essere infinite, quindi trovare il valore massimo di queste è un problema...

Tipologie di policy

Le policy possono essere:

- 1) dipendenti dalla storia si chiama -> HD (history dependent)
- 2) dipendenti dal tempo e non dalla storia, si chiama -> MARKOV
- 3) se non dipende da nulla si dice "stazionaria" e si indica con pigreco (π)

Quindi nonostante l'insieme delle policy è un insieme infinito e quindi non è possibile trovare il massimo, allora come faccio? Bisogna quindi prendere un sottoinsieme di policy e in particolare quelle stazionarie e deterministiche, ma in che modo?

Possiamo farlo se agiamo su stati e azioni **finiti**. Le policy stazionarie e deterministiche sono tutte le possibili azioni associate agli stati. (che si rappresenta con A elevato alla S , dove A sono le azioni e le S sono gli stati)

Definition

Any policy obtaining optimal state-value or optimal action-value is called a *optimal policy*: π_* is a optimal policy if

$$q_{\pi_*} = q_* \quad \text{or} \quad v_{\pi_*} = v_*$$

Se esiste una policy "ottimale", cioè il valore di tutti gli stati è il valore V^*/Q^* (massimo) di quello stato.

La policy ottimale è deterministica in quanto sceglie sempre l'azione (una) che massimizza il ritorno, come si trova? si trova agendo "greedy" sulle azioni.

quindi:

Optimal policy

Theorem

For any Markov decision process:

- There exists (at least) an optimal policy π_* .
- All optimal policies π_* achieve the optimal state-value and the optimal action-value function: $v_{\pi_*} = v_*$ and $q_{\pi_*} = q_*$.

Equazione di ottimalità per funzione stato valore ottimale V^* di Bellman

A questo punto troviamo l'equazione ottimale di Bellman per la policy π^*

Partiamo dalla "classica" equazione di Bellman che agisce su tutte le policy, ovvero:

$$v_{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{\pi}(s') \right)$$

Modifichiamola facendola agire solo sulla policy ottimale:

$V^* =$

$$v_{\pi_*}(s) = \sum_{a \in \mathcal{A}} \pi_*(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{\pi_*}(s') \right)$$

ma la policy valore ottimale sceglie le azioni in maniera greedy, quindi:

$$v_*(s) = \max_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \right)$$

$V^* =$

da notare che la V^* è ricorsiva in quanto compare all'interno della formula.

Equazione di ottimalità per funzione azione-valore ottimale Q^* di Bellman

$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a' \in \mathcal{A}} (q_*(s', a'))$$

Ricapitolando, le equazioni di Bellman ottimali e non sono:

Bellman equations: a recap

Bellman equation for v_π

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') \right)$$

Bellman optimality equation for v_*

$$v_*(s) = \max_a \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \right)$$

Bellman equation for q_π

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a'|s') (q_\pi(s', a'))$$

Bellman optimality equation for q_*

$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a'} (q_*(s', a'))$$

Risoluzione dell'equazione di Bellman

Partiamo dal presupposto che l'MDP è un processo finito, ciò significa che gli stati sono definiti e numerati. I metodi per risolvere l'equazione sono di tipo "iterativo".

NB: In questa fase non stiamo facendo ancora apprendimento, stiamo facendo quella che viene chiamata "pianificazione" in quanto conosciamo il modello. Per implementare la pianificazione, vengono utilizzate tecniche di programmazione dinamica.

La programmazione dinamica è una tecnica di risoluzione dei problemi, si applica in due fasi:

- il problema complesso viene scomposto in sotto-problemi più semplici
- i sotto-problemi vengono poi ricomposti per risolvere il problema originale

Ci sono però dei requisiti, ovvero:

- il problema deve poter essere scomposto, si dice in questo caso che deve avere "una sotto-struttura ottimale"
- i sotto-problemi devono essere sovrapponibili, ovvero alcuni calcoli si possono ripetere in quanto l'algoritmo riceve in input gli stessi valori, per cui in questi è possibile implementare delle ottimizzazioni per evitare l'utilizzo di computazionale, come per es. l'utilizzo di cache, o altro.

Per calcolare la funzione valore di un policy dobbiamo quindi applicare l'equazione di Bellman iterativamente

Programmazione dinamica parte 2 (fase di pianificazione non apprendimento)

Stiamo quindi cercando di risolvere la fantomatica equazione di Bellman per determinare il valore dello stato e/o dello stato ottimale. La tipologia di casi che stiamo affrontando vengono detti "tabellari" in quanto, avendo a che fare con stati finiti, possono essere rappresentati in una tabella con valori finiti.

Per casi più complicati si utilizzano i casi "non tabellari" utilizzando gli approssimatori, ovvero le reti neurali.

Veniamo al pseudo algoritmo che utilizza l'equazione di Bellman.

Iterative policy evaluation for estimating $V \sim v_\pi$

Input: Policy π to be evaluated.

Parameter: Threshold $\theta > 0$ determining accuracy of estimation.

Output: Estimate V of v_π .

Initialize $V(s)$, for all $s \in \mathcal{S}$, arbitrarily except $V(\text{final}) = 0$.

do

$\Delta \leftarrow 0$

for $s \in \mathcal{S}$ **do**

$v \leftarrow V(s)$

$V_{\text{temp}}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a V(s'))$

$\Delta \leftarrow \max(\Delta, |v - V_{\text{temp}}(s)|)$

end

$V \leftarrow V_{\text{temp}};$

while $\Delta > \theta$

spiegone:

L'algoritmo inizializza il vettore V con dei valori casuali tranne l'ultimo stato che deve essere inizializzato a zero, Questo perchè l'algoritmo è di tipo "bootstrap" ovvero utilizza i valori del vettore al tempo t_1 per determinare i valori degli stati al tempo t_2 . Lo stato "finale" deve essere **assorbente**, ovvero dallo stato finale o si esce o fa ripartire il loop di convergenza.

Ma questa versione non si usa, di seguito la versione migliorata, che in pratica utilizza l'ultimo valore per aggiornare V in quanto l'assegnazione di V avviene nel ciclo più interno, vedi sotto:

Input: Policy π to be evaluated.

Parameter: *Threshold* $\theta > 0$ determining accuracy of estimation.

Output: Estimate V of v_π .

Initialize $V(s)$, for all $s \in \mathcal{S}$, arbitrarily except $V(\text{final}) = 0$.

do

$\Delta \leftarrow 0$

for $s \in \mathcal{S}$ **do**

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a V(s'))$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

end

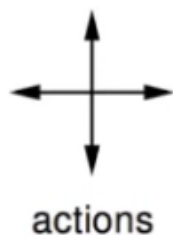
while $\Delta > \theta$

in questo modo tutti gli stati successivi al primo beneficiano di valori aggiornati degli stati precedenti.

Esercizio:

Di seguito abbiamo un "mondo griglia" con due stati finali.

Policy evaluation example: gridworld



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

- 14 nonterminal states $1, \dots, 14$, one terminal state (shown twice as \square), four actions $\rightarrow, \leftarrow, \uparrow, \downarrow$.
- Actions leading out of the grid leave state unchanged.
- Goal: reach a terminal state as soon as possible.

Exercise

Model the gridworld task as an Episodic MDP. Hint: \mathcal{S} , \mathcal{A} and dynamics are given (describe them), we need to add rewards.

Analisi:

gli stati sono 15, le azioni sono 4, un esempio di policy potrebbe essere: al 70% vado giù, e all'30% vado a destra.

Invece un esempio di policy ottimale:

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

Vedendo il modello possiamo pianificare l'azione, per questo siamo in modalità "pianificazione".

Ora proviamo con la policy: al 70% vado giù, e all'30% vado a destra.

Applicazione pratica dell'equazione di Bellman valore-stato (con policy non ottimale)

$$\begin{aligned} \forall s \quad \pi(\downarrow | s) &= 0.7 \\ \pi(\rightarrow | s) &= 0.3 \\ v_1(1) &= \sum_a \pi(a | 1) \left(R_1^a + \sum_{s' \in S} P_{1s'}^a v_0(s') \right) \\ &= \sum_a \pi(a | 1) \left(-1 + \sum_{s'} P_{1s'}^a \cdot (0) \right) = \\ &= 0.7(-1+0) + 0.3(-1+0) \\ &= -1 \end{aligned}$$

Policy Evaluation

Ora veniamo al secondo passaggio, che va anche a spiegare il primo (rappresentato qui sopra)

Il secondo passaggio (v_2) dopo che il primo passaggio ha settato il valore V degli stati a -1 . (ad eccezione ovviamente degli stati finali T che valgono zero)

Il passaggio logico è:

La prima parte dell'espressione è relativa alla policy, che nel caso specifico è al 70% andare giù e al 30% andare a destra. (si legge come probabilità di raggiungere lo stato successivo s' facendo l'azione in s , dove s in questo caso vale 1 in quanto è lo stato 1)

$$\sum_a \pi(a|1)$$

Quindi per prima cosa inseriamo nell'espressione le due probabilità 0.7 e 0.3 che moltiplicano una certa quantità.

$$\begin{aligned} v_2(1) &= \sum_a \pi(a|1) \left(R_1^a + \sum_{s'} P_{1s'}^a v_2(s') \right) \\ &= 0.7 \cdot \left(\right) \pi \\ &\quad + 0.3 \cdot \left(\right) \end{aligned}$$

Diagrammatic annotations: An arrow points from the action 'a' in the first term to a downward arrow labeled 'a = ↓'. Another arrow points from the action 'a' in the second term to a rightward arrow labeled 'a = →'.

(questa quantità in formula si legge come la ricompensa media dell'andare nello stato s' sommato alla probabilità dell'azione "a" di andare nello stato s' moltiplicata (la probabilità) per la reward dello stato s' . Attenzione che la probabilità di andare nello stato s' da s , in questo particolare caso, è 100% quindi 1. Da cui ne deriva che:

$$v_2(1) = \sum_a \pi(a|1) \left(R_{11}^a + \sum_{s'} P_{1s'}^a v_1(s') \right)$$

$$= 0.7 \cdot \left(-1 + \underset{\substack{\downarrow \\ 1}}{\overset{\downarrow}{P_{15}^a}} v_1(5) \right) \pi$$

$a = \downarrow$

$$+ 0.3 \cdot \left(-1 + 1(-1) \right)$$

$$= -2$$

Il tutto va iterato per tutti gli stati dell'ambiente fino a che i valori degli stati "convergono", il che significa che variano di poco.

Il metodo consente di propagare in modalità "backward" dagli stati finali (quelli più vicini allo stato T) verso quegli iniziali i valori.

Alla fine ciascuno stato del modello avrà un valore che indicherà quanti "passi" sono necessari per arrivare alla meta, dandoci quindi un'indicazione sul comportametto della policy.

Conclusioni dell'esercizio

Con una policy che assegna il 25% per ogni azione gli stati valori ottenuti sono:

- [[0. -14. -20. -22.]
- [-14. -18. -20. -20.]
- [-20. -20. -18. -14.]
- [-22. -20. -14. 0.]

che indicano la ricompensa totale media che si otterrebbe partendo dallo stato. per es. se parto dallo stato "1" (quello vicino all'uscita) otterrei mediamente, con la policy al 25% una ricompensa media di circa -14.

ATTENZIONE! stiamo parlando di una ricompensa media in quanto ci potrebbero essere casi in cui arrivo subito all'uscita andando per es. subito a sx oppure casi in cui inizio a va girare per la griglia, esistono quindi **traiettorie** lunghissime.

Ecco l'esempio di codice:

allego a questa pagina le classi di ambiente utilizzate per simulare il "mondo griglia"

```
from envs.gridworld import GridworldEnv
import numpy as np

# 0 su 1 dx 2 giù 3 sx
# instanzio l'ambiente griglia
env = GridworldEnv()

print (env.nS)
print (env.nA)
env.reset()
env._render()
env.step(1)
print()
env._render()
env.step(1)

def policy_evaluation(policy, env, discount_factor=1.0, theta=0.00000001):
    """
    :param policy: [A,S] matrice di shape (rango = 29 a due dimensioni che rappresenta la
    policy, dove sulle righe ci sono gli stati e sulle colonne le azioni
    :param env: rappresentazione dell'ambiente
    :param discount_factor: fattore di sconto gamma
    :param theta: valore che termina la valutazione della policy una volta che la funzione
    stato valore è sotto questa soglia
    :return:
    """
```

```

# inizializzamo la funzione stato valore V, per semplificare la vita li portiamo a zero
# dove il valore è il valore di OGNI STATO
V = np.zeros(env.nS)

# il ciclo deve fermarsi solo quanto la qta delta calcolata è minore o uguale a theta
passaggi = 0
while True:
    delta = 0
    passaggi += 1

    # per ciascuno stato effettuo un "full backup"
    # questo primo ciclo fa passare tutti gli stati
    for s in range(env.nS):
        v = 0

        # questo il ciclo fa la sommatoria delle azioni
        # ovvero esegue tutte le azioni possibili nello stato s
        for a, action_prob in enumerate(policy[s]):

            # per ciascuna azione chiedo all'ambiente in quale stato finisco, la
            # probabilità di finirci, la ricompensa e se eventualmente è lo stato finale
            # NB: la situazione NON è vera in quanto stiamo facendo "esperienza" senza
            # fare un passo, questo non è vero RL è pianificazione
            for prob, next_state, reward, done in env.P[s][a]:

                # equazione di Bellman
                # probabilità dell'azione per la probabilità della transizione
                v += action_prob * prob * (reward + discount_factor * V[next_state])

        # quanto del valore stato funzione è variato tra tutti gli stati
        delta = max(delta, np.abs(v - V[s]))
        V[s] = v

    # finiamo la valutazione una volta che il valore è sotto la soglia theta
    if delta < theta:
        break

return np.array(V), passaggi

```

```

# definiamo una policy generica, in pratica per ogni stato del mondo griglia associa 4
probabilità di eseguire una azione, le probabilità sono tutte al 25%
random_policy = np.ones( [env.nS, env.nA])/env.nA

v,passaggi = policy_evaluation(random_policy, env)

v = v.reshape(4,4)
print(v,passaggi)

```

Migliorare la policy (policy improvement-iteration)

Il metodo per migliorare la policy è quello di compiere tutte le azioni possibili e scegliere quella che restituisce il valore maggiore agendo in maniera "greedy", ovvero:

Improve the policy by acting greedily with respect to v_π :

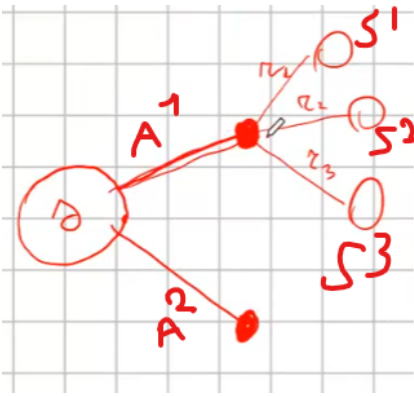
$$\forall s \in \mathcal{S} \quad \pi'(s) := \operatorname{argmax}_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') \right)$$

Come faccio a calcolare π' (primo)? be faccio tutte le azioni "a" possibili, faccio "**one step looking forward**", ovvero faccio un passo avanti compiendo un'azione sola, guardo cosa succede al passo successivo chiedendo all'ambiente cosa può succedere con quale probabilità. Facendo queste azioni prendo quella che massimizza il ritorno.

Ma quando una policy è migliore di un'altra? beh quando per ogni stato s , applicando la nuova policy π' , il valore dello stato è maggiore-uguale al valore della vecchia policy. Nella pratica faccio N iterazioni compiendo N traiettorie fino a quando i valori convergano, a questo punto confronto tutti gli stati valori con quelli precedenti, se sono migliori allora faccio un altro giro, loppando fino a quando gli stati valori sono uguali al precedente, il che identifica la policy migliore.

Ma cosa vuol dire nella pratica?

Supponiamo di avere uno stato " s " e un'azione " a_1 " che porta in tre stati (s', s_2, s_3), dove per ognuno dei tre stati avviamo una probabilità di entrare in ciascuno stato con una relativa reward, abbiamo anche il valor medio totale dello stato π di ciascun stato che era stato precedentemente calcolato applicando la policy. (vedi schema sotto riportato)



Ora per ogni azione (es. a_1 e a_2) voglio calcolare il valore medio dell'azione stessa che significa calcolare il valore delle "cose/fatti" che possono accadere pesati per le probabilità. Nell'esempio sopra riportato significa prendere la probabilità che accada s' con la relativa reward r dato lo stato s con l'azione a_1 , ovvero:

$p(s',r | s, a_1) * (r + V\pi(s'))$ -> probabilità di andare in s' con l'azione a_1 che moltiplica la somma della reward ottenuta per andare nello stato s' + il valore medio totale dello stato s' .

Questa operazione va ripetuta per ogni stato raggiunto dall'azione a_1 , ovvero nel nostro caso s', s_2, s_3 e sommata per questi tre stati quindi:

$$Q\pi(s,a) = \sum (xx) p(s_{xx},r | s, a_{xx}) * (r + V\pi(s_{xx})) \text{ dove } xx \text{ è } s', s_2, s_3$$

Poi analogamente calcoliamo il $Q\pi(s,a_2)$ e ne calcoliamo il massimo tra $argmax (Q\pi(s,a), Q\pi(s, a_2))$ e scegliamo **l'azione** associata al **$Q\pi$ massimo**. In questo modo andiamo a migliorare la policy in maniera "**greedy**" ovvero cercando di massimizzare il valore. Il nome di questo algoritmo di chiama "policy iteration". Il cui pseudo codice è sotto riportato:

Parameter: Threshold $\theta > 0$. **Output:** Estimate of v_* and π_* .

1. Initialize $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}$ arbitrarily, $V(\text{final}) = 0$.
2. Policy evaluation.
3. Policy improvement.

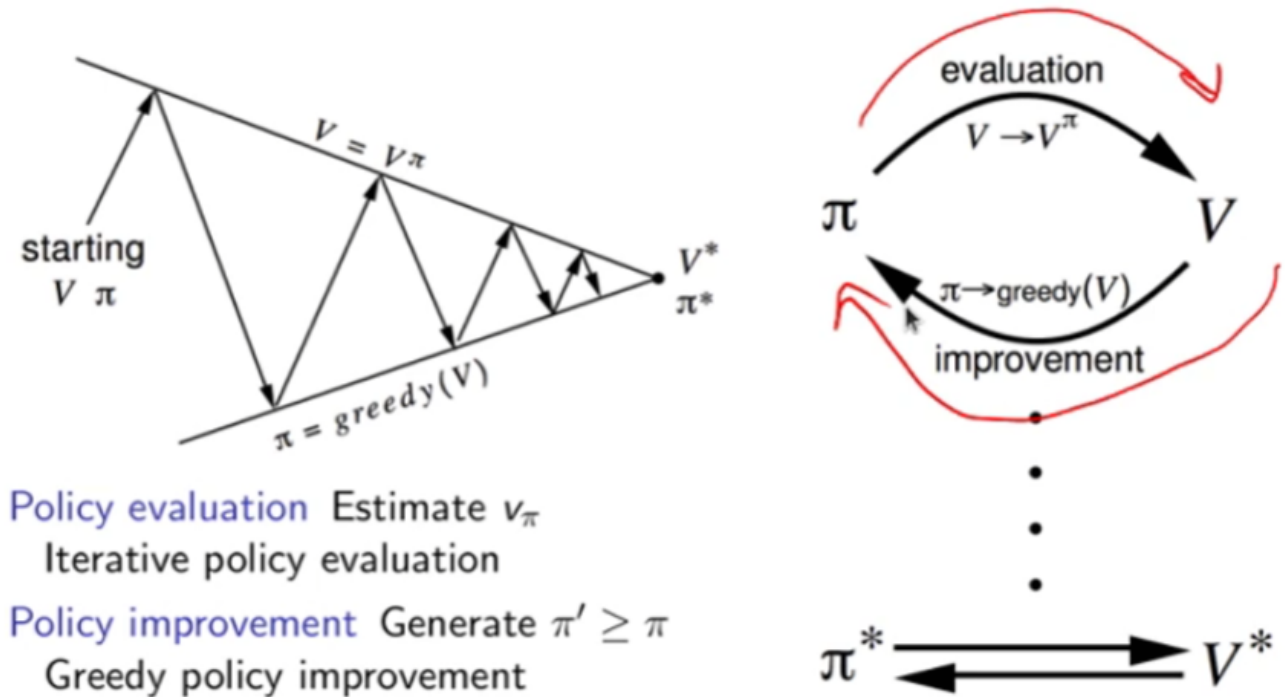
```

policy-stable ← true
for  $s \in \mathcal{S}$  do
  | old-action ←  $\pi(s)$ 
  |  $\pi(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a V(s'))$ 
  | if  $\text{old-action} \neq \pi(s)$  then
  | | policy-stable ← false
  | end
end
if policy-stable = true then
  | return  $V, \pi$ 
else
  | go to 2
end

```

L' algoritmo alterna la valutazione (evaluation) della policy al improvement fino a che non ottengo la policy.

Policy iteration



dove:

1) inizializza i valori randomicamente e seleziona una policy casuale

2) applica la valutazione delle policy, ovvero calcola V_π

3) miglioramento, ovvero:

- facciamo un ciclo sugli stati S (detta anche "sweep") -> NB essendo pianificazione di un MDP a stati finiti rende la cosa fattibile, in realtà in giochi come "go" non è possibile in quanto gli stati sono un numero enorme, ma questa è un'altra storia)
- ciclo su tutte le azioni dello stato che sto valutando
 - salvo la vecchia azione ($V(s_{xx})$)
 - effettuo la nuova azione e vedo quale valore $V(s_{xx})$ ottengo
 - confronto vecchia e "nuova funzione valore" $V(S)$ se il valore della nuova azione è migliore allora l'algoritmo non ha ancora trovato la policy migliore e quindi aggiorno la policy con l'azione migliore e continuo a loopare

- fatti passare tutti gli stati se la policy non è stabile riparto dal punto 2
- in generale il tutto loopa finchè i valori non convergono il che significa che tra uno sweep e l'altro non ci sono grandi variazioni

Ricordo che la pianificazione si può applicare solo quando disponiamo del modello, che è il caso più facile e poco probabile che accada nel RL.

Di seguito l'algorithmo completo

```

from envs.gridworld import GridworldEnv
import numpy as np

def policy_evaluation(policy, env, discount_factor=1.0, theta=0.001):
    """
    :param policy: [A,S] matrice di shape (rango = 29 a due dimensioni che rappresenta la
    policy, dove sulle righe ci sono gli
        stati e sulle colonne le azioni
    :param env: rappresentazione dell'ambiente
    :param discount_factor: fattore di sconto gamma
    :param theta: valore che termina la valutazione della policy una volta che la funzione
    stato valore è sotto questa soglia
    :return: ritorna la tabella stato valore (V) contenente i valori ottimali per arrivare
    allo più efficacemente allo stato terminale
        in pratica l'algorithmo incentiva ad andare verso lo stato il cui valore è più
    alto, es. tra due stati che valgono
        rispettivamente- 14 e -20 si va verso il valore più grande ovvero -14.

    NOTA: - le azioni sono: 0 = su, 1 = dx, 2 = giù, 3 sx
        - la reward è sempre -1 tranne negli stati terminali
        - negli stati terminali le azioni non effettuano spostamento di stato.
    """

    # inizializzamo la funzione stato valore V, per semplificare la vita li portiamo a zero
    # dove il valore è il valore di OGNI STATO
    V = np.zeros(env.nS)

    # il ciclo deve fermarsi solo quanto la qta delta calcolata è minore o uguale a theta
    passaggi = 0
    while True:

```

```

delta = 0
passaggi += 1

# per ciascuno stato effettuo un "full backup"
# questo primo ciclo fa passare tutti gli stati
for stato in range(env.nS):
    v_appo = 0

    # questo il ciclo fa la sommatoria delle azioni
    # ovvero esegue tutte le azioni possibili nello stato s
    azioni_stato = enumerate(policy[stato])
    for azione , policy_prob in azioni_stato:

        # per ciascuna azione chiedo all'ambiente in quale stato finisco, la
        # probabilità di finirci, la ricompensa
        # e se eventualmente è lo stato finale
        # NB: la situazione NON è vera in quanto stiamo facendo "esperienza" senza
        # fare un passo, questo non è
        # vero RL è pianificazione
        # Da notare che questo ciclo è utile solo nel caso in cui a fronte di una
        # azione ci possono essere diversi
        # stati destinazione con diverse probabilità, nel nostro caso ci sarà sempre e
        # solo uno stato destinazione
        # con probabilità (env_prob) pari a 1
        for env_prob, next_state, reward, done in env.P[stato][azione]:

            # equazione di Bellman
            # probabilità dell'azione per la probabilità della transizione
            v_appo = v_appo + policy_prob * env_prob * (reward + discount_factor *
V[next_state])

        # delta misura l'errore su tutti gli stati, ovvero misura la differenza tra tutti
        # gli
        # nella passata precedente N-1 e l'attuale N.
        # Nella pratica significa che per ogni stato il delta massimo trovato è quello
        # indicato nella variabile
        # e quindi potrà poi essere comparato con theta per terminare la valutazione della
        # policy
        delta = max(delta, np.abs(v_appo - V[stato]))

```

```

V[stato] = v_appo

# stampo la tabella degli stati
# vr = V.reshape(env.shape).round(3)
# print(vr,stato,passaggi,'%5.15f' % delta)
# input()

# finiamo la valutazione una volta che il valore è sotto la soglia theta
if delta < theta:
    break

return np.array(V), passaggi

def policy_improvement (env, policy, policy_evaluation_fn, discount_factor=1.0):
    """
    Funzione che migliora la policy iterativamente

    :param env: ambiente openAI
    :param random_policy: passiamo una policy iniziale
    :param policy_evaluation_fn: funzione che valuta la policy e restituisce gli stati
    funzione valore V
    :param discount_factor: fattore di sconto gamma
    :return: ritorna una tupla che che tiene la nuova policy ottimale, gli stati valori e il
    totale passaggi effettuati
    """
    def one_step_lookahead(state, V):
        """
        funzione helper che calcola il valore di tutte le azioni dato uno stato

        :param state:
        :param V:
        :return: Q-value -> un vettore di lunghezza env.nA ovvero tutte le azioni possibili
        nello stato che contiene il valore funzione stato per ogni azione.
        """
        A = np.zeros(env.nA)
        # prendo tutte le azioni possibili nello stato passato in input alla funzione
        for a in range(env.nA):

            # come reagisce l'ambiente in quello stato

```

```

    # per ogni azione ottengo il valore totale medio
    # NB: interrogo l'ambiente e sarà lui a dirmi con quella azione in quali stati
andrò e con quale probabilità e reward
    for prob, next_state, reward, done in env.P[state][a]:
        # calcolo il valore medio di tutti gli stati raggiungibili effettuando
l'azione "a".
        A[a] += prob * (reward + discount_factor * V[next_state])
    return A

passaggi = 0
# ciclo di controllo.
while True:

    passaggi += 1

    # valutiamo la policy
    V, psg = policy_evaluation_fn(policy, env, discount_factor)

    # setto il flag che mi interromperà il loop while
    policy_stable = True

    # estraggo tutti gli stati
    # SWEEP
    for s in range(env.nS):

        # effettuo lo "one step ahead" per capire quali valori mi restituisce l'azione che
ho deciso di
        action_values = one_step_lookahead(s, V)

        ##### calcolo l'azione migliore
        # IMPROVMENT
        #####
        best_a = np.argmax(action_values)

        # scelgo l'azione con probabilità maggiore tra quelle presente nella VECCHIA
(attuale) policy
        # per quello stato, nella pratica estrae l'indice
        prev_policy_action = np.argmax(policy[s])

        # confronto l'azione della vecchia policy con quella trovata dall'improvement e

```

```

verifico se coincidono
    # se non coincidono allora la policy non è ancora stabile e quindi non ottimale.
    if prev_policy_action != best_a:
        policy_stable = False

    # setto l'indice dell'azione che massimizza il valore calcolato a 1 per cambiare
la policy ottimizzandola
    policy[s] = np.zeros(env.nA)
    policy[s][best_a] = 1

    # se la policy è stabile allora abbiamo trovato quella ottimale, quindi interrompo il
ciclo principale ed esco
    if policy_stable:
        break

return policy, V, passaggi

# 0 su 1 dx 2 giù 3 sx
# instanzio l'ambiente griglia
env = GridworldEnv([4,4])

print (env.nS)
print (env.nA)
env.reset()
env._render()
env.step(1)
print()
env._render()
env.step(1)

# definiamo una policy generica, in pratica per ogni stato del mondo griglia associo 4
probabilità di eseguire una azione, le probabilità sono tutte al 25%
random_policy = np.ones( [env.nS, env.nA])/env.nA
print(random_policy)

policy_ottimale, V, passaggi = policy_improvement (env, random_policy, policy_evaluation)
#V,passaggi = policy_evaluation(random_policy, env)

```

```
v = V.reshape(env.shape).round(0)
print(v,passaggi)
#
print()
print(policy_ottimale)
```

(<https://www.youtube.com/watch?v=QY3yxYyK4wM&list=PLMee1hSjLKdAL16E-7EzqHXsGOgzo8iro&index=13>)

Iterazione di valore

Con queste tipologie di algoritmi vogliamo **ottimizzare** le policy combinando la policy evaluation e la policy improvement.

Nel iterazione della policy ricordiamo che l'algoritmo è diviso in due parti, la prima che riguarda la valutazione della policy (evaluation) e la seconda che la migliora agendo in maniera "greedy" ovvero scegliendo l'azione che massimizza il valore medio atteso prendendo la argmax dei valori restituiti dal "one step look ahead".

Nella iterazione di valore invece l'algoritmo di semplice e indirettamente si ottimizza combinando i due step secondo lo pseudo codice sotto riportato:

Parameter: Threshold $\theta > 0$ determining accuracy of estimation.

Initialize: Initialize $V(s)$, for all $s \in \mathcal{S}$, arbitrarily except $V(\text{terminal}) = 0$.

Output: Estimate of v_* .

do

$\Delta \leftarrow 0$

for $s \in \mathcal{S}$ **do**

$v \leftarrow V(s)$

$V(s) \leftarrow \max_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a V(s'))$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

end

while $\Delta > \theta$

che in pratica applica l'equazione di Bellman sostituendo direttamente il valore medio atteso della funzione stato valore.:

Bellman optimal update

Use the Bellman optimality equation for a **single-step update**:

$$v_{k+1}(s) = \max_a \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

Then:

$$v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_*$$

di seguito l'algoritmo rivisto, la differenza in termini di performance rispetto al precedente è notevole.

```
from envs.gridworld import GridworldEnv
import numpy as np

import time

# funzione per misurare il tempo di esecuzione di una funzione tramite un decorator
# basterà:
# 1) importare la funzione -> from objs.TimerDecorator import timing_decorator
# 2) decorare la funzione da misurare con @timing_decorator
def timing_decorator(func):
    def wrapper(*args, **kwargs):
        start = time.time()
        original_return_val = func(*args, **kwargs)
```

```

    end = time.time()
    print("time elapsed in ", func.__name__, ": ", end - start, sep='')
    return original_return_val

return wrapper

from envs.gridworld import GridworldEnv
import numpy as np

from utils.TimerDecorator import timing_decorator

@timing_decorator
def value_iteration(env, theta = 0.00001, discount_factor=1.0):
    """
    Funzione che migliora la policy policy_improvement

    :param env: ambiente openAI
    :param theta: valore che termina la valutazione della policy una volta che la funzione
    stato valore è sotto questa soglia
    :param discount_factor: fattore di sconto gamma
    :return: ritorna una tupla che che tiene la nuova policy ottimale, gli stati valori e il
    totale passaggi effettuati
    """

    def one_step_lookahead(state, V):
        """
        funzione helper che calcola il q-value (valore) di tutte le azioni dato uno stato

        :param state:
        :param V:
        :return: Q-value -> un vettore di lunghezza env.nA ovvero tutte le azioni possibili
        nello stato che contiene il valore funzione stato per ogni azione.
        """
        A = np.zeros(env.nA)
        # prendo tutte le azioni possibili nello stato passato in input alla funzione
        for a in range(env.nA):

            # come reagisce l'ambiente in quello stato
            # per ogni azione ottengo il valore totale medio

```

```

        # NB: interrogo l'ambiente e sarà lui a dirmi con quella azione in quali stati
andrò e con quale probabilità e reward
        for prob, next_state, reward, done in env.P[state][a]:
            # calcolo il valore medio di tutti gli stati raggiungibili effettuando
l'azione "a".
            # equazione di "punto fisso"
            A[a] += prob * (reward + discount_factor * V[next_state])
        return A

# inizializziamo gli stati con dei valori casuali a piacere (array di N stati)
V = np.zeros(env.nS)
while True:
    # inizializzo la variabile che condiziona lo stop del loop
    delta = 0

    # aggiornare di stati
    for s in range(env.nS):
        # guardo uno step avanti restituendo la funzione stato valore di ogni azioni che
può essere fatta nello stato
        best_action_value = np.max(one_step_lookahead(s,V))

        # confronto l'azione con il valore dello stato nell'iterazione precedente
        delta = max(delta, np.abs (best_action_value - V[s]))

        # salvo lo stato valore
        V[s] = best_action_value

    # se il delta è sotto theta allora i valori sono prossimi alla convergenza desiderata
    if delta < theta:
        break

# ora dagli stati valore ricavo la policy  $V^*(s) \rightarrow Q^*(s,a)$ 
# inizializzo la policy
policy = np.zeros([env.nS,env.nA])
for s in range (env.nS):
    azione = np.argmax(one_step_lookahead(s,V))
    policy[s][azione] = 1.0

return policy, V

```

```

# 0 su 1 dx 2 giù 3 sx
# instanzio l'ambiente griglia
env = GridworldEnv([4, 4])

# print(env.nS)
# print(env.nA)
env.reset()
# env._render()
# env.step(1)
# print()

policy_ottimale, V, = value_iteration(env)
# V,passaggi = policy_evaluation(random_policy, env)

print(V)
#
# print()
print(policy_ottimale)

```

Esercizio 2

Abbiamo un certo capitale inferiore a 100 (es. 99 o 1 o 44) e lo scopo è arrivare esattamente a 100. (valori oltre il 100 non sono considerati, il goal è esattamente a 100) Possiamo scommettere solo il capitale a nostra disposizione lanciando una moneta scommettendo quello che vogliamo sulla base di quello che abbiamo. (es. ho 70 lancio la moneta e scommetto per es. 30 per arrivare a 100, se perdo vado a 40 se vinco vado a 100 e ho terminato il gioco vincendo. Ovviamente se arrivo a zero ho perso.

Il problema può essere modellato come un MDP, senza 1) fattore di sconto, 2) episodico, ovvero da traiettorie (successione di azioni, rewards e stati) di lunghezza finita che si concludono in uno (o più) stato/i terminale/i raggiungibile/i e 3) finito ovvero che stati e azioni sono finiti.

In generale quando le ricompense sono tutte zero, tranne che nella transizione dallo stato precedente allo stato terminale e lo stato terminale stesso, dove in questo caso vale 1, allora la funzione stato valore indica la probabilità di vittoria in quello stato.

Ricodiamo che: il valore dello stato è media con la probabilità di tutte le traiettorie possibili della ricompensa totale.

```
import numpy as np
```

```
"""
```

Abbiamo un certo capitale inferiore a 100 (es. 99 o 1 o 44) e lo scopo è arrivare esattamente a 100.

(valori oltre il 100 non sono considerati, il goal è esattamente a 100) Possiamo scommettere solo

il capitale a nostra disposizione lanciando una moneta scommettendo quello che vogliamo sulla

base di quello che abbiamo. (es. ho 70 lancio la moneta e scommetto per es. 30 per arrivare a 100, se perdo vado a 40 se vinco vado a 100 e ho terminato il gioco vincendo.

Ovviamente se arrivo a zero ho perso.

Il problema può essere modellato come un MDP, senza 1) fattore di sconto, 2) episodico, ovvero da traiettorie (successione di azioni, rewards e stati) di lunghezza finita che si concludono in uno (o più) stato/i terminale/i raggiungibile/i e 3) finito ovvero che stati e azioni sono finiti.

In pratica lo stato rappresenta il capitale posseduto, e in base a quello la policy dovrebbe cercare di assumere il comportamento ottimale.

In prima battuta va definito il modello. Il modello restituisce sempre ricompensa 0 tranne quando arrivo nello stato terminale 100. Notare che la ricompensa non è quanto stiamo guadagnando in quanto l'obiettivo è arrivare a 100 e non accumulare più soldi possibili.

```
"""
```

```
def values_iteration_for_gamblers(probabil_moneta = 0.5 , theta=0.0001, discount_factor=1.0):
```

```
    """
```

```
        :param probabil_moneta: probabilità che esca testa
```

```
        :param theta: valore delta di convergenza
```

```
        :param discount_factor: fattore di sconto
```

```

:return: policy e funzione stato valore
"""

def one_step_look_ahead(in_capital, V, rewards, probab_monetina, discount_factor):
    """
    In base allo stato in cui mi trovo, ovvero il capitale a mia disposizione, effettuo
una possibile
    giocata per tutte le possibili giocate a mia disposizione. Ovvero se ho un capitale di
40, farò
    una giocata partendo dalla somma 1 fino al massimo a mia dispisizione. (ovvero 40)

    La funziona ci dice nella pratica per ogni azione fatta con il capitale passato quanto
vale quella azione.

    :param capitale: capitale del giocatore (rappresentato dallo stato)
    :param V: stati valore
    :param rewars: ricompense

    :return: ritorna l'elenco dei valori medi ottenibili effettuato tutte le azioni nello
stato ovvero
        del capitale passato in input
    """

    # inizializzo le azioni possibili in ogni stato
    A = np.zeros(101)

    # passato il capitale a disposizione, calcolo tutte le possibili giocate da quella
minima (1) al massimo
    # che posso giocare con il capitale a disposizione passato in input.
    # il min(capital, (100-capital)) serve per calcolare il capitale giocabile
considerando che
    # non posso andare oltre 100 anche se il capitale a mia disposizione lo consetirebbe
    capitale_giocabile = min(in_capital, (100 - in_capital))

    # effettuo tutte le giocate possibili
    for giocata_in_soldi in range(1, capitale_giocabile+1):
        # a questo punto per ogni giocata vedo il valore medeio di quello che può
succedere effettuando
        # tutte le giocate possibili.

```

```

        A[giocata_in_soldi] =      probab_monetina * (rewards[in_capital +
giocata_in_soldi] + V[in_capital + giocata_in_soldi] * pow(discount_factor, giocata_in_soldi))
\
        + (1 - probab_monetina) * (rewards[in_capital -
giocata_in_soldi] + V[in_capital - giocata_in_soldi] * pow(discount_factor, giocata_in_soldi))

        # array che indica per ogni azione fatta con il capitale a mia disposizione, quanto
vale ciascuna azione.
        return A

#####
##### INIZIO #####
#####

# creo un array di 101 elementi che rappresenta le rewards dove:
# - l'elemento zero è lo stato di perdita del gioco -> stato terminale
# - l'elemento 101 è lo stato di vincita del gioco -> stato terminale
# tutti gli stati avranno una rewards pari a zero tranne l'ultimo (100) che avrà reward =
1
rewards = np.zeros(101)

# ricordo che è zero based, quindi l'elemento 100 è il 101
rewards[100] = 1

# inizializzo la funzione stato valore con dei valori a "caso" facciamo zero per facilitare
V = np.zeros(101)

# creo la policy ottimale
policy = np.zeros(101)

# ciclo principale
while True:

        # inizializzo la variabile che verrà utilizzata per interrompere i ciclo di
valutazione della policy
        # fino alla sua convergenza.
        delta = 0

        # simulo tutti le possibili giocate con tutti i possibili capitali a mia disposizione
        # fa passare tutti gli stati

```

```

for giocata in range (1,100):

    # verifico quanto valgono tutte le giocate con il capitale simulato
    # faccio tutte le azioni possibili
    A = one_step_look_ahead(giocata, V, rewards, probab_monetina, discount_factor)

    # determino il valore della giocata più alto
    best_action_value = np.max(A)

    # verifico se la giocata è migliore della precedente
    delta = max(delta, np.abs(best_action_value-V[giocata]))

    # salvo il valore migliore nella funzione stato valore
    V[giocata] = best_action_value

# verifico se interrompere il loop perchè la funzione stato valore sta converendo
if delta < theta:
    break

# questo ciclo indica come giocare in ogni stato, dove lo stato rappresenta il capitale
posseduto
# serve solo per identificare l'azione migliore da intraprendere basandosi sulle stati
valori ottimali precedentemente
# calcolatir
for giocata in range (1,100):
    A = one_step_look_ahead(giocata, V, rewards, probab_monetina, discount_factor)
    best_action_value = np.argmax(A)
    policy [giocata] = best_action_value

return policy,V

policy, V = values_iteration_for_gamblers(discount_factor=1)

# tampo la funzione valore
for i in range(101):
    if i % 10 == 0:
        print()
    print (" %i-%s"%(i,V[i]), end='')

print ()

```

```
# stampo la policy
for i in range(101):
    if i % 10 == 0:
        print()
        #print (" %i-%s"%(i,policy[i]), end='')
        print(" %i-%s" % (i,policy[i]), end='')
# print (policy)
# print (np.reshape(policy,(10,10)))
```

Revision #91

Created 2023-06-27 18:21:21 UTC by marco

Updated 2025-05-04 07:21:15 UTC by marco