# Assignment 3

## Introduzione al problema

- Guidare una macchina da corsa per fare una curva
- L'obiettivo è andare più veloci possibili senza uscire dal tracciato
- L'auto si trova in una delle posizioni discrete della griglia.
- La velocità è discreta: numero di celle della griglia spostate orizzontalmente e verticalmente per ogni intervallo di tempo.
- Le azioni sono incrementi delle componenti di velocità.
- Le velocità sono limitate a un massimo di 5.
- Gli episodi iniziano in modo casuale sulla linea di partenza.
- Gli episodi terminano se si taglia il traguardo.
- Le ricompense sono -1 per ogni intervallo.
- Se si raggiunge il limite, si torna casualmente sulla linea di partenza.

#### Creazione del tracciato

- Viene identificato tramite una matrice t.c. fuori=-1, start=1, tracciato= 0, fine=2
- Alla funzione si passa le dimensioni del tracciato
- La linea di partenza è composta di 5 celle
- Per ogni riga successiva, tramite vettori di probabilità, si decide se il tracciato si restringe, allarga o rimane lo stesso
- Vincolo di lasciare il percorso di almeno 4 celle
- Il traguardo è nelle ultime 5 righe a destra

#### Creazione del tracciato

```
def create_track(N1,N2): 3 usages  sqronz
track = np.zeros( shape: (N1,N2), dtype=int) #fuori =-1, start=
s1=randint( a: 0, N2-10)
s2=randint(s1+5,N2)
    if i in range(s1,s2):
for i in range(1,N1-5): #lato destro lascio le ultime 5 che
    if track[i,y-5]==-1:
    track[N1-i,N2-1]=2
for riga in track:
    print(riga)
```

### On-policy every-visit Monte Carlo control

- Si inizializza la policy casualmente
- Si utilizza una policy ε-greedy, in cui l'epsilon decade a ogni episodio fino a un valore minimo
- Si genera l'episodio
- Si calcolano i ritorni

$$\begin{split} N(S_t, A_t) &\leftarrow N(S_t, A_t) + 1 \\ Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} \left(G - Q(S_t, A_t)\right) \\ A^* &\leftarrow \arg\max_a Q(S_t, a) \\ \text{for all } a &\in \mathcal{A}(S_t) \text{ do} \\ \pi(a|S_t) &\leftarrow \begin{cases} 1 - \varepsilon + \frac{\varepsilon}{|\mathcal{A}(S_t)|}, & \text{if } a = A^*, \\ \frac{\varepsilon}{|\mathcal{A}(S_t)|}, & \text{otherwise} \end{cases} \end{split}$$

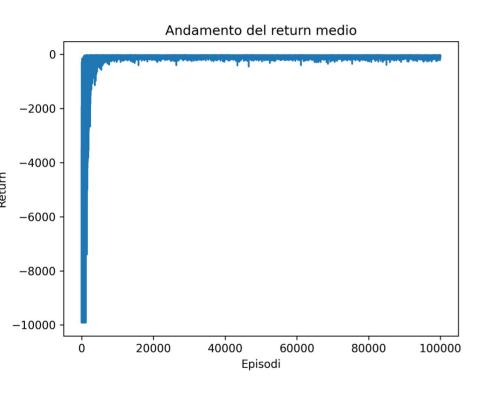
### On-policy every-visit Monte Carlo control

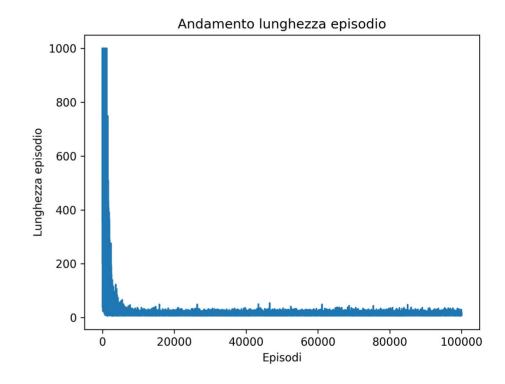
```
Q = np.zeros((N1,N2,A)) #(N1,N2) griglia, azioni
N = np.zeros((N1, N2, A))
pi = np.zeros( shape: (N1, N2), dtype=int)
actions = [(ao, av) for ao in [-1, 0, 1] for av in [-1, 0, 1]]
avg_lengths = []
eps = 1
eps end = 0.01
for el in range(100_000):#numero episodi
   episode = generate_episode(track,pi,eps)
   total return = 0
   for t, (stato, azione , reward) in enumerate(reversed(episode)):
       s1, s2 = stato
       a_idx = actions.index(azione)
       total return += G
       N[s1, s2, a idx] += 1
       alpha = 1 / N[s1, s2, a_idx]
       Q[s1, s2, a_idx] += alpha * (G - Q[s1, s2, a_idx])
       best_a = np.argmax(Q[s1,s2,:])
       pi[s1, s2] = best_a
    avg_returns.append(total_return)
    avg_lengths.append(len(episode))
   eps = max(eps_end, eps * decay_rate)
   if (el+1) % 1000 == 0:
             f"Lunghezza media: {np.mean(avg_lengths[-50:]):.1f}, "
return pi, Q, avg_returns, avg_lengths
```

### Generazione episodio

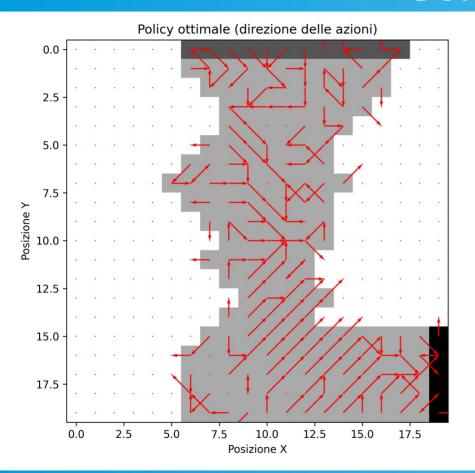
- Si parte dalla linea di partenza con velocità 0
- Limite numero di step poiché le prime policy casuali potrebbero non trovare soluzione
- Si sceglie l'azione secondo la policy ε-greedy
- Si calcola la nuova velocità limitandola a 5
- La macchina si muove prima verticalmente e poi orizzontalmente
- Se esce fuori riparte dalla linea di partenza con velocità 0
- 2 tipologie di reward

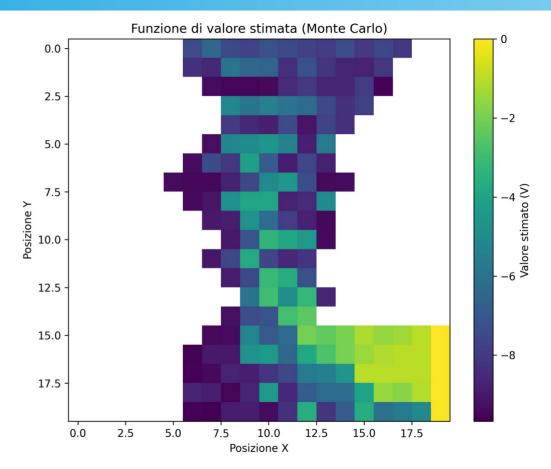
## Andamento training



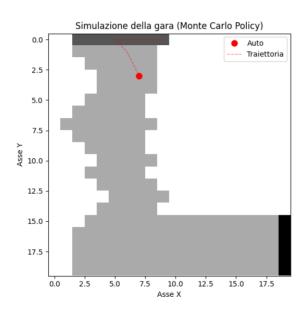


### Soluzione





#### Corsa



#### FINE