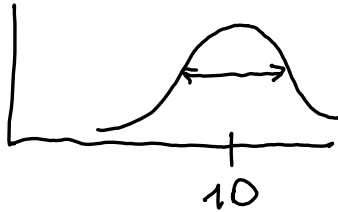


M.A.B

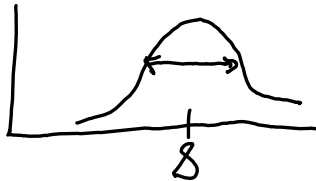
Restaurant example

STANDARD DEV



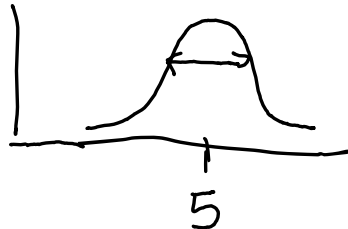
$$S_x = 5$$

$$\bar{X} = 10$$



$$S_x = 4$$

$$\bar{X} = 8$$



$$S_x = 2.5$$

$$\bar{X} = 5$$



Queste info solo
nascoste

X e' il REWARD DEI RISTORANTI
in termini di SODDISFAZIONE

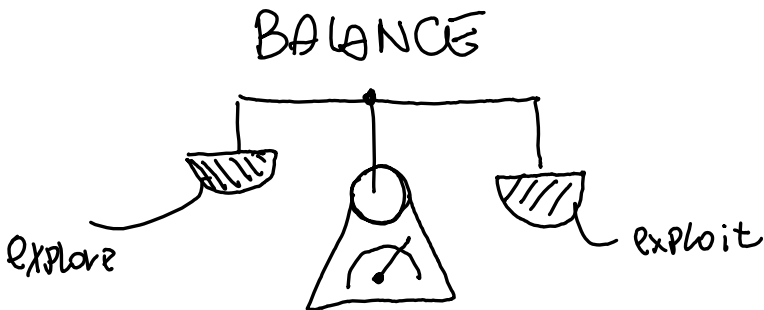
NOI ANDIAMO IN QUESTI RISTORANTI, SENZA
SAPERE LA LORO CAPACITÀ DI REWARD

① EXPLORATION

- Andiamo prima e così e poi
Secondo una logica

② EXPLOITATION

- UNA volta che copiamo che quel ristorante
che abbiamo esplorato più volte ci
soddisfa ci CONTINUIAMO AD ANDARE PER
massimizzare il reward.



IN BASE ALLA SCELTA DI RISOLUZIONE DEL PROBLEMA AVREMO UN DISCOSTAMENTO DALLA SOLUZIONE IDEALE (S.I.)

" S.I. = SAPENDO GIÀ LE INFO EXPLOITATO SOLO SUL MIGLIORE. "

QUESTA METRICA DI ERRORE LE INDICHEREMO CON
 $P(R\hat{\theta})$

ϵ - GREEDY STRATEGY

- SE ESPLORIAMO E BASTA, AD UN CERTO PUNTO SAPPIAMO QUALE BANDIT CI DA UN REWARD MIGLIORE, QUINDI PRENDIAMO TEMPO CONTINUANDO AD ESPLORE
- SE EXPLOITIAMO SU UNO NON PRENDIAMO ABBASTANZA DATI PER SAPERE QUALE POTREBBE ESSERE IL MIGLIORE

CI SERVE COSTRUIRE UNA STRATEGIA DI
BILANCIAMENTO.

- PRENDIAMO UNA ϵ
SOLITAMENTE BASSA

$$5\% - 10\% - \frac{1}{\# \text{ TRY}} \%$$

ϵ è la % PER CUI SCEGLIAMO RANDOM
VN RISTORANTE

$1 - \epsilon$ è la probabilità di continuare
l'exploitation sul reward migliore.

ES : 300 TRY

$$\left. \begin{array}{l} 30 \text{ EXPLORE} \\ 270 \text{ EXPLOIT} \end{array} \right\} \epsilon 10\%$$

UCB1 Alg / method

still:

300 DAYS

3 RESTAURANTS \rightarrow HAPPINESS REWARD
distributions

the problem with ϵ -Greedy

SE CI BASIAMO SUL VALORE MEDIO DI
HAPPINESS, POTREMO AVER VISITATO
POCHE VOLTE UN RISTORANTE E QUINDI
QUELLA MEDIA POTREBBE NON ESSERE VALIDA.

" ES: se R_1 è stato visitato 3
volte e $\mu_{R_1} = 2.9$ //

questo 2.9 non è selezionato
da un n° suff di visite //

QUI ENTRA QUINDI UNA NUOVA STRATEGIA
DI BILANCIAMENTO.

- AD OGNI TRY (t) , PRENDIAMO UN RISTORANTE (R) , tale che :

$$\left[\hat{\mu}_R + \sqrt{\frac{2 \ell_m(t)}{N_t(r)}} \right] \text{ E' MASSIMO}$$

→ HOPPING'S
INEQUALITY

LA MEDIA
DEL NEWBARD
PER R_i
BY GIVEN
TRY

QUESTA PARTE
del polinomio
e' l'incertezza
che abbiamo
basata sul n°
di visite

- $N_t(r)$
N° volte che
abbiamo visitato R_i

- $\ell_m(t)$
NATURAL log of try

LA SECONDA PARTE DEL POLINOMIO
E' DETTA HOPFING'S INEQUALITY

- CI FORNISCE UN UPPER BOUND DI
PROBABILITA' CHE LA SOMMA DELLE
VARIABILI CASUALI INDIPENDENTI SI DISCOSTI
DAL SUO VALORE ATTESO

MA TORNIAMO A NOI - - -

SE ABBIAMO 2 RISTORANTI VISITATI
CON LA STESSA MEDIA / NON CI BASIAMO
SOLO SU DI ESSA MA IL SECONDO FATTORE
DEL POLINOMIO SARA' DETERMINANTE

N.B

QUEL VALORE E' DA MASSIMIZZARE

IN TERMINI DI EXPLORATION E EXPLOITATION

$\hat{\mu}_R$: EXPLOITATION HELPER

$\sqrt{\frac{2 \ln(t)}{N_t(r)}}$: EXPLORATION HELPER

QUANDO VISITIAMO UN RISTORANTE UNA
VOLTA IN +, AUMENTIAMO IL DENOMINATORE
IL CHE DETERMINA CHE MENO E MENO
OBIETTIVO- SEMPRE MENO LIKELY VISITARLO
PERCHÉ LA FRAZIONE DIMINUISCE.
