

Decision Theory

9/4/25

- Techniques to take decisions in the face of uncertainty.
- when faced with uncertainty, don't use ad hoc techniques.
 - Decision theory provides well-developed techniques.
 - But, what is uncertainty, & why do we need to reason with it?

Which route to take?

$$R_1: S \rightarrow F \rightarrow B : 310$$

$$R_2: S \rightarrow R \rightarrow P \rightarrow B : 278$$

- Ideally would take shortest route.
- What if traffic is a factor now?
- Suppose you have accurate info that
 - R_1 has fast moving traffic while R_2 has slow moving traffic.
- A logical agent may reason
 - slow(R_2), fast(R_1)
 - slow(x) \Rightarrow Avoid(x)
 - Avoid(x) ^ fast(y) \Rightarrow select(y)
 - Agent selects R_1 .
- Do you think this is a good way to reason?
- Why or why not?

Ex: 80% 4 hrs 20% 10 hrs 5.2 hrs

70% 4.5 hrs 30% 5 hrs < 5 hrs

80%

Uncertainty

- May not have categorical answers.
 - May not know exactly which route is slow.
 - 50% chance that one route is slow
 - Even if we know it is slow, how slow?
- Uncertainty changes the way an agent makes a decision.
- Probability theory used in helping agents reason with uncertainty.
 - Summarize uncertainty due to our ignorance or laziness.
 - May be very very costly to remove such uncertainty.
- Rational decision depends on
 - Relative importance of various goals
 - Likelihood or degree to which they will be achieved.

How to make decisions?

Decision theory = Probability theory +
(deals with chance)

utility theory

(deals with outcomes).

Fundamental idea's

- The ~~MEU~~ Maximum expected utility principle.

• Agent is rational if & only if it chooses the action that yields the highest expected utility averages over all possible outcomes of action.

• weigh the utility of each outcomes by the probability of that it occurs.

Utility theory - von Neumann (Game theory)

Let X be the set of outcomes. \succ be the preference of a player over the set of outcomes.

Axioms:

- Completeness: every pair of outcomes is ranked.
- Transitivity: If $x_1 \succ x_2$ and $x_2 \succ x_3$ then $x_1 \succ x_3$
- Substitutability: If $x_1 \succ x_2$ then any lottery in which x_1 is substituted by x_2 is equally preferred.
- Decomposability: Two different lotteries assign same probability to each outcome, then player is indifferent b/w these two lotteries.
- Monotonicity: If $x_1 \succ x_2$ and $p > q$, then $[x_1:p, x_2:1-p] \succ [x_1:q, x_2:1-q]$
- Continuity: If $x_1 \succ x_2 \succ x_3$, $\exists p \ni x_2 \sim [x_1:p, x_3:1-p]$

Von Neumann and Morgenstern:

Theorem: Given a set of outcomes X and a preference relation on X that satisfies

above six axioms, there exists a utility function $u: X \rightarrow [0, 1]$ with the following properties:

$$u(x_1) > u(x_2) \text{ iff } x_1 \succ x_2.$$

$$U([x_1; p_1, x_2; p_2; \dots; x_m; p_m]) = \sum_{j=1}^m p_j u(x_j)$$
$$\boxed{\sum p_j = 1}$$
$$= \sum_{j=1}^m p_j u(x_j)$$

Basic Probability Notation

- Suppose we have prob. over all possible values (sample space) of a variable denoted by $P(X)$.

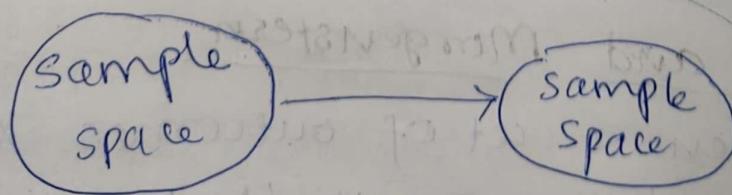
- Called the prob. distribution for the random variable:

- $P(\text{Weather} = \text{sunny}) = 0.6$
- $P(\text{Weather} = \text{rain}) = 0.1$
- $P(\text{Weather} = \text{cloudy}) = 0.29$
- $P(\text{Weather} = \text{snow}) = 0.01$
- $P(\text{Weather}) = \langle 0.6, 0.1, 0.29, 0.01 \rangle$

$$\sum P(X = x_i) = 1$$

$x_i \in \{\text{sunny, rain, cloudy, snowy}\}$

Random Variable



→ This is not event space.

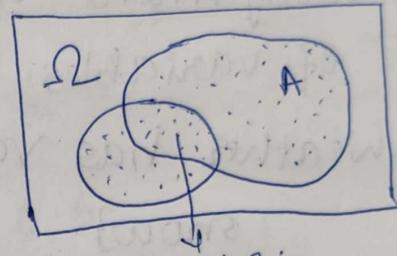
(Ω, \mathcal{F}, P)
 2 event space (discrete)
 2 $[0, 1]$ -continuous - Borel sigma algebra

Conditional Probability

$$A, B \\ P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{if } P(B) > 0$$

$k_1 - (A \cap B)$

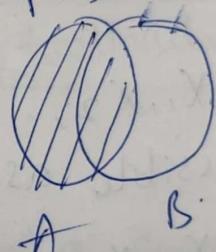
$k - B$
n - total sample space



$$\frac{k_1}{k} = \frac{k_1/n}{k/n} = \frac{P(A \cap B)}{P(B)} \quad A \cap B = A \cap B$$

Product Rule:

- $P(A|B) = P(A \cap B) / P(B)$ if $P(B) > 0$ definition
- $P(A \cap B) = P(A|B) * P(B)$
- $P(A \cap B) = P(B|A) * P(A)$



Posterior or conditional probability

- $P(A|B)$ probability of A given all we know is B.
- $P(X = \text{sunny} | \text{last 2 days were sunny})$.
- If we know B and also know C, then
 $P(A|B \cap C)$

Diagnosis domain: Introducing Cavity variable

- Doctor dia - photo
- Not everytime

Joint Probability Distribution

- Notation for distribution on multiple variables.
- P(Weather, cavity) i.e., $P(\text{weather} \wedge \text{cavity})$
- Joint prob. distribution is a table.
- Assign probabilities to all possible assignment of values for combinations of variables.
- Weather has values {sunny, rain, cloudy, snow}.
- Cavity has values {true, false}.
- 4×2 table of probabilities called joint probability distribution.
- $P(X_1, X_2, \dots, X_n)$ assigns probabilities to all possible assignment of values to variables X_1, X_2, \dots, X_n .

Inference using full joint distribution

		toothache		¬ toothache			
		catch	¬ catch	catch	¬ catch	catch	¬ catch
		0.108	0.012	0.072	0.008	0.016	0.064
Cavity	¬ Cavity	0.012	0.002	0.008	0.000	0.000	0.000
¬ Cavity	Cavity	0.016	0.064	0.144	0.576	0.000	0.000

3 variable
cavity,
Toothache,
catch

Dentist

$$\Pr(X=x) = \sum_{(y,z)} \Pr(X=x, Y=y, Z=z)$$

probe (catch)
or not

Marginalization

$2 \times 2 \times 2$
table

$P(\text{cavity}) = 0.108 + 0.012 + 0.072 + 0.008 = 0.2$ (called marginal probability of cavity).

$$P(\text{cavity} \mid \text{toothache}) = \frac{P(C \cap T)}{P(T)} = \frac{0.12}{P(T)}.$$

$$\begin{aligned} P(T) &= 0.108 + 0.012 + 0.016 + 0.064 \\ &= \frac{0.12}{0.2} = \frac{12}{20} = 0.6 \end{aligned}$$

$$P(\text{cavity} \mid \text{toothache}) = 0.6$$

$$P(\neg \text{cavity} \mid \text{toothache}) = 0.4 [1 - 0.6]$$

$$= \frac{0.016 + 0.064}{P(T)} = \frac{8}{20} = 0.4$$

$$\boxed{P(\text{cavity} \mid \text{toothache}) + P(\neg \text{cavity} \mid \text{toothache}) = 1}$$

Marginalization

- Denominator can be viewed as a normalization constant $1 / P(\text{toothache})$.

$$\begin{aligned} P(\text{cavity} \mid \text{Toothache}) &= \alpha P(\text{cavity}, \text{toothache}) = \alpha [P(\text{cavity}, \text{toothache}, \text{catch}) + \\ &\quad P(\text{cavity}, \text{toothache}, \neg \text{catch})] = \\ &\quad 0.108 + 0.012 = \alpha 0.12 \end{aligned}$$

$$\begin{aligned} \alpha P(\neg \text{cavity}, \text{toothache}) &= \alpha [0.016 + 0.064] \\ &= \alpha 0.08 \end{aligned}$$

- Normalizing $\alpha [0.12, 0.08]$ gives $\langle 0.6, 0.4 \rangle$

- Issue: For a domain with n Boolean variables, i/p table size of $O(2^n)$.
full joint distribution in tabular form not a practical tool for building reasoning systems.

Independence:

- Helps in reducing size of domain representation and complexity of inference problem.
- Let's add weather into the example:
sunny, rain, cloudy, snow $\rightarrow 2^* 2^* 2^* 4$ entries in table.
- $P(\text{toothache, catch, cavity} | \text{cloudy}) = P(\text{cloudy} | \text{toothache, catch, cavity}) \cdot P(\text{toothache, catch, cavity}).$
- $P(\text{cloudy} | \text{toothache, catch, cavity}) = P(\text{cloudy})$
- 32 element table becomes 8 element + 4 element table.
- $$P(A|B) = \frac{P(A)P(B)}{P(B)} = P(A).$$

Bayes' Rule

- Given that
 - $P(A \wedge B) = P(A|B) * P(B)$
 - $P(A \wedge B) = P(B|A) * P(A)$

$P(B|A) = \frac{P(A \cap B)}{P(A)}$

We can determine $P(B|A)$ given $P(A|B)$,
 $P(B)$ and $P(A)$.

- $P(\text{effect}|\text{cause})$ may be causal knowledge,
 $P(\text{cause}|\text{effect})$ diagnostic knowledge.
- Often have conditional probabilities on
causal relationships $\rightarrow P(\text{symptoms}/\text{disease})$
- Need to infer diagnostic knowledge
many times.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)}$$

Probability theory

$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)}$$

Posterior \nwarrow Prior \downarrow Evidence

Diagnostic knowledge \nearrow Likelihood \downarrow Causal knowledge.

Ex:

$$P(I|F) = \frac{P(F|I) \cdot P(I)}{P(F)}$$

Ex: S: Proposition that patient has stiff neck
 M: Proposition that patient has meningitis
 • Meningitis causes stiff-neck, 70% of the time

Given,

$$P(S|M) = 0.7$$

$$P(M) = 1/50,000$$

$$P(S) = 0.01$$

$$P(M|S) = \frac{P(S|M) \cdot P(M)}{P(S)} = \frac{0.7 \times \frac{1}{50,000}}{0.01}$$

$$= \frac{7 \times 10^{-1} \times \frac{1}{5} \times 10^{-4}}{10^{-2}} = 7 \times 10^{-1} \times 0.2 \times 10^{-2}$$

$$= 1.4 \times 10^{-3} = 0.0014$$

→ Why not estimate $P(M|S)$ the same way as we do $P(S|M)$?
 ↓ because pr. of meningitis is low

↳ Why do we need Baye's rule?

- $P(S|M)$ is causal knowledge, does not change
 - It is "model based"
 - It reflects the way meningitis works
- $P(M|S)$ is diagnostic, tells us likelihood of M given symptom S.
 - Diagnostic knowledge may change with circumstance, thus helpful to derive it.
 - If there is an epidemic, probability of Meningitis goes up.

rather than trying to estimate $P(M|S)$ which is very hard, we can compute it using the other terms.

- However, $P(S|M)$ is unaffected by epidemic.
 $P(E|C) \rightarrow$ causal lang.
because of this,
this is effect(E)

Computing the denominator : $P(S)$

We wish to avoid computing the deno. in the Bayes' rule

- May be hard to obtain
- In our example, the deno. is $P(S)$
- Techniques to compute { or avoid computing $P(S)$ }.

Our first approach is to compute relative likelihoods

If M (meningitis) and W (whiplash) are two possible explanations:

$$P(M|S) = P(S|M) * P(M) / P(S)$$

$$P(W|S) = P(S|W) * P(W) / P(S)$$

$$\therefore P(M|S) / P(W|S) = \frac{P(S|M) * P(M)}{P(S|W) * P(W)}$$

Another approach : Approach #2

$$P(M|S) = P(S|M) * P(M) / P(S)$$

$$P(\text{NOT}(M)|S) = P(S|\text{NOT}(M)) * P(\text{NOT}(M)) / P(S)$$

- $P(M|S) + P(\text{Not}(M)|S) = 1$ [These 2 quantities must sum to 1].
- $[P(S|M) * P(M) / P(S)] + [P(S|\text{Not}(M)) * P(\text{Not}(M)) / P(S)] = 1$
- $P(S|M) * P(M) + P(S|\text{Not}(M)) * P(\text{Not}(M)) = P(S)$
- calculate $P(S)$ in this way?

Find $P(B|A_i) = P(A_i|B) \cdot P(B)$

$$\sum_{i=1}^k P(A_i|B) \cdot P(B_i)$$

$$U.A_i = \Omega, \quad A_i \cap A_j = \emptyset.$$

- $P(M|S) + P(W|S) + P(D|S) + P(ND|S) = 1$.

Simple Example:

- Suppose 2 identical urns
- URN1 colored red from inside, has $\frac{1}{3}$ black balls, $\frac{2}{3}$ red balls
- URN2 colored black from inside has $\frac{2}{3}$ red balls, $\frac{1}{3}$ black balls.
- We select one URN at random; can't tell how it is colored inside.
- What is the probability that URN is colored inside? 0.5
- What if we were to select a ball at random from URN, and it is red? Does that change the probability?

$$P(\text{Red-um} | \text{Red-ball}) = \frac{P(\text{Red-ball} | \text{Red-um}) * P(\text{Red-um})}{P(\text{Red-ball})}$$

$$= \frac{\frac{2}{3} * 0.5}{P(\text{red-ball})}$$

→ How to calculate $P(\text{red-ball})$?

$$P(\text{Black-um} | \text{Red-ball}) = P(\text{Red-ball} | \text{Black-um}) * \frac{P(\text{Black-um})}{P(\text{Red-ball})}$$

$$= \frac{\frac{1}{3} * 0.5}{P(\text{red-ball})}$$

Thus by our approach #2: $\frac{2}{3} * 0.5 / P(\text{red-ball}) + \frac{1}{3} * 0.5 / P(\text{red-ball}) = 1$

Thus, $P(\text{red-ball}) = 0.5$

$\therefore P(\text{Red-um} | \text{red-ball}) = 2/3$

More General forms of Bayes' rule:

- $P(Y|X) = P(X|Y) * P(Y) / P(X)$

- Bayes' rule for multi-valued variables

- Generalize to some background

evidence e

- $P(Y|X, e) = P(X|Y, e) * P(Y|e) / P(X|e)$

Conditional Independence

- Toothache & Cough are independent given the presence of cavity.

- Each is directly caused by cavity, toothache is caused by cavity

and a density catches using probe,

- A toothache cannot be caught by a probe nor a probe results in a toothache.

- Conditional independence captured as:

$$P(\text{toothache} \wedge \text{catch} | \text{cavity}) = P(\text{toothache} | \text{catch, cavity}) \cdot P(\text{catch} | \text{cavity})$$

$$\therefore P(\text{toothache} | \text{catch, cavity}) = P(A | B, C)$$

- If A & C are conditionally independent given B

- Then, probability of A is not dependent on C.

$$\therefore P(A | B \wedge C) = P(A | B)$$

- Therefore, $P(\text{toothache}, \text{catch} | \text{cavity}) =$

$$P(\text{toothache} | \text{cavity}) \cdot P(\text{catch} | \text{cavity})$$

A = toothache, B = cavity, C = catch

Combining Evidence

- S : Proposition that patient has stiff neck

- H : Proposition that patient has severe headache.

- M : Proposition that patient has meningitis

- Meningitis causes stiff-neck, 50% of the time.

- Meningitis causes head ache, 70% of the time.

- Probability of meningitis should go up if both symptoms reported - how to

$$P(A|B \wedge C) = P(A \wedge B \wedge C) / P(B \wedge C)$$

Numerator:

$$P(M \wedge S \wedge H) = P(S|M \wedge H) * P(M \wedge H)$$

$$= P(S|M) * P(M \wedge H) \rightarrow \text{why}$$

$$= P(S|M) * P(H|M) * P(M)$$

Going back to our example:

$$P(M|S \wedge H) = P(S|M) * P(H|M) * P(M) / P(S \wedge H)$$

$S \wedge H$ are conditionally independent

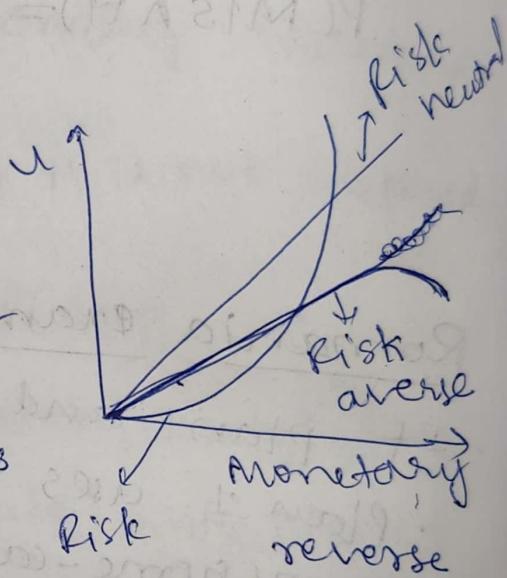
Romania Example:

- If plan1 and plan2 are the two plans
- Plan 1 uses route 1.
- $P(\text{home-early} | \text{plan 1}) = 0.8$,
 $P(\text{stuck 1} | \text{plan 1}) = 0.2$
- Route 1 will be quick if flowing, but stuck for + hour if slow.
- $V(\text{home-early}) = 100$,
 $V(\text{stuck 1}) = -1000$
- Assigned numerical values to outcomes
- Plan 2 uses route 2.
- $P(\text{home-somewhat-early} | \text{plan 2}) = 0.7$,
 $P(\text{stuck 2} | \text{plan 2}) = 0.3$
- Route 2 will be somewhat quick if flowing, but not bad even if slow.
 $V(\text{home-somewhat-early}) = 50$,
 $V(\text{stuck 2}) = -10$

- Appn of MEV principle
- $EV(\text{Plan 1}) = P(\text{home early}|\text{Plan 1}) * V(\text{home early})$
 $= P(\text{stuck 1}|\text{Plan 1}) * V(\text{stuck 1})$
 $= 0.8 * 100 + 0.2 * (-1000) = -120$
 - $EV(\text{Plan 2}) = P(\text{home-somewhat-early}|\text{Plan 2}) * V(\text{home-somewhat-early}) +$
 $P(\text{stuck 2}|\text{Plan 2}) * V(\text{stuck 2})$
 $= 0.7 * 50 + 0.3 * (-10) = 32$

Risk Aversion

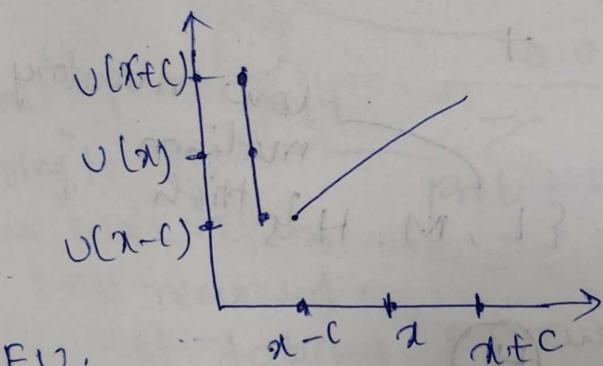
- We are risk averse
- Our utility fns. are money as follows:
 - Our first million means a lot $f: U(\$1M) = 10$
 - Second million not so much $U(\$2M) = 15$ (not 20)
 - Third million even less so $U(\$3M) = 18$ (not 30)
- Additional money is not buying us as much utility.
- If we plot amount of money on the x-axis and utility on the y-axis, we get a concave curve



More Risk Aversion

- Key: Slope of utility fn is continuously decreasing.
- We will refuse to play a monetarily fair bet.

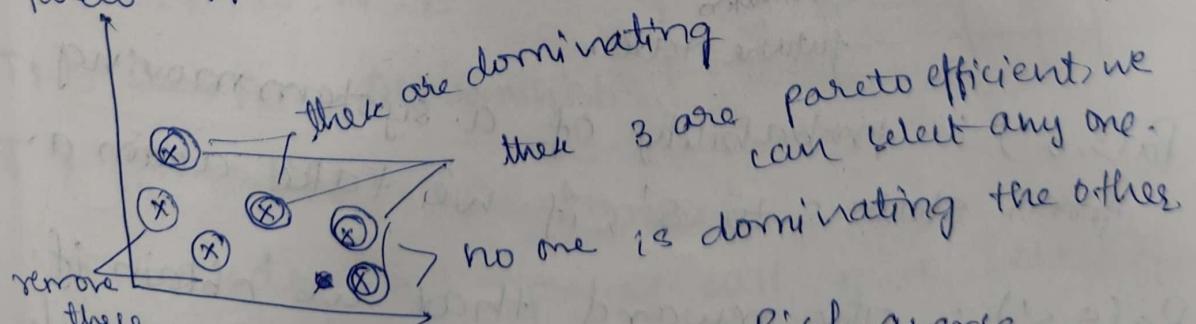
- suppose we start with x dollars.
- We are offered a game:
 - 0.5 chance to win 1000 dollars ($c=1000$)
 - 0.5 chance to lose 1000 dollars ($c=-1000$)
 - Expected monetary gain or loss is zero (hence monetarily fair).
 - Should be neutral to it, but seems we are not! why?
 - $V(x+c) = V(x) + V(x) - V(x-c)$
 - $V(x+c) + V(x-c) < 2V(x)$
 - $[V(x+c) + V(x-c)]/2 < V(x)$
 - $EU(\text{Playing the game}) < EU(\text{not playing the game}).$



MEU:

16/4/25

- Rule out dominated values.
- The one which maximizes the weighted sum.
- Pareto efficient $\stackrel{\text{ex:}}{=}$ two features multiple options



Utility theory

Risk averse
Risk Neutral
Risk seeking/loving

Pareto efficient:

Suppose we are dealing with N -dimensional utilities. Option a is pareto efficient if

$$\nexists y \text{ s.t. } u_i(y) > u_i(x) \quad i=1, \dots, N$$

with strict inequality for at least one;

→ Pareto efficient is one method to select.

→ $f(u_1, \dots, u_N)$ - mapping. Need not be linear. (One method of solution).

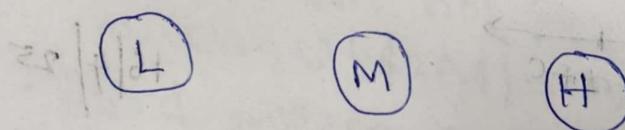
Actions:

Ex:

	Production	Marketing
a_0	+	
a_1	+	+
a_2	-	+
a_3	-	-

Markov decision process:

Let states be, $S = \{L, M, H\}$



$$P: S \times S \times A \rightarrow [0, 1]$$

$$R: \text{current future} \times S \times S \times A \rightarrow R$$

$P(a(s, s'))$: probability of a system moving from state s to s' , if we take action a .

$R(a(s, s'))$: The reward that we obtain if we take an action a in state s & reach state s' .

$\langle S, A, R, P \rangle$ - Markov decision process

- depends on the previous step only, not on all previous states.

- P_A - markovian process R_A - not markovian, just showing a change in system, it is not a process

Example: Photo

N, E, S, W

(1, 3), (1, 2), (1, 1)

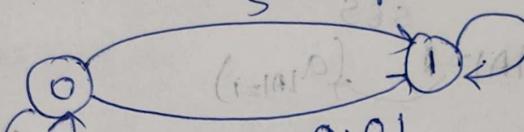
(2, 3), (2, 1)

(3, 3), (3, 2), (3, 1)

will MEV work here?

s - reward

Ex:



- but after

going to ①, probability

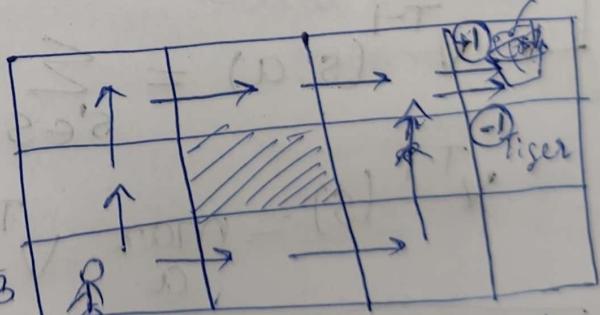
we can't return

no reward & staying in 0 is so, not taking s (reward) is better

we can assure MEV is the best.

Red: N, N, E, E, E
E, E, N, N, E

$\pi: S \rightarrow A$ Deterministic



$\pi: S \rightarrow D(A)$ Randomized.

No randomized policy will give better than optimal deterministic function

+ rounds (finite horizon)
infinite horizon

~~Π_0, Π_1, Π_2~~ possible values - $\langle q, s, A, \pi \rangle$

$$\Pi^0 = \Pi^1 = \Pi^2 \rightarrow \text{(stationary policy).}$$

$|A|^1$ are possible policies.

$$\begin{aligned}\Pi_0 &= (0, a_0) (1, a_1) = \Pi_2 = (0, a_1) (1, a_0) \\ \Pi_1 &= (0, a_1) (1, a_0) \quad \Pi_3 = (0, a_0) (1, a_0)\end{aligned}$$

Finite Horizon:

Each action gives some rewards (Expected reward)

$$\begin{aligned}
 &\text{① } \sum_{s \in S} P_{a_0}(1, s) \times R_{a_0}(1, s) \\
 &\text{② } \sum_{s \in S} P_{a_1}(1, s) \times R_{a_1}(1, s) \\
 &\vdots \\
 &\text{⑬ } \sum_{s \in S} P_{a_{|A|-1}}(1, s) \times R_{a_{|A|-1}}(1, s) \\
 &\text{⑭ } V^T(s) = \max_{a \in A} \sum_{s' \in S} P_a(s, s') R_a(s, s') \\
 &\text{⑮ } \Pi^T(s) = \arg \max_{a \in A} V^T(s, a)
 \end{aligned}$$

value of states

$$V^{T-1}(s, a) = \sum_{s' \in S} P_a(s, s') [V^T(s') + R_a(s, s')]$$

$$V^{T-1}(s) = \max_a V^{T-1}(s, a)$$

$$V^{T-k}(s) = \max_a \left[\sum_{s' \in S} P_a(s, s') \left\{ V^{T-k+1}(s') + R_a(s, s') \right\} \right]$$

$\rightarrow T$ is given.

$$\textcircled{1} \quad \begin{aligned} a_0, a_1 & P_{a_0}(0,1) = 0.6 = P_{a_1}(1,0) \\ \pi(0) = a_0 & P_{a_0}(0,0) = 0.4 = P_{a_1}(1,1) \\ \pi(1) = a_1 & \end{aligned}$$

$$R_{a_0}(0,1) = R_{a_0}(0,0)$$

$$R_{a_1}(1,0) = R_{a_1}(1,1) = 1$$

let $V(0) = 2, V(1) = 3$ (and consider $\pi(0) = a_0, \pi(1) = a_1$, these values depend upon policy.)

$$V^T(0) = 0.4 [1 + \frac{V(0)}{2}] + 0.6 [1 + \frac{V(1)}{3}]$$

$$(V^T(0)) = 3.6$$

$$V^T(1) = 0.4 [1 + \frac{V(1)}{2}] + 0.6 [1 + \frac{V(0)}{3}]$$

$$R_{a_0}(1,0)$$

$$R_{a_1}(1,1)$$

$\rightarrow 5, 5, 5, 5, 0, 5, 5, 5, 5, 0, \dots$ } are both

$\rightarrow 4, 4, 4, 4, 4, 4, \dots$ } equal or the rewards diff?

$$\begin{aligned} & \rightarrow 0, 0, 0, 0, 0, 5, 5, 5, 5, 0, \dots \\ & 4, 4, 4, 4, 4, 4, \dots \end{aligned}$$

$\gamma = 0$ (ignoring future value)

$\gamma \rightarrow$ discount factor for the values.

generally b/w $[0 \& 1]$.

$$V^T(0) = 0.4 [1 + \gamma \cdot 2] + 0.6 [1 + \gamma \cdot 3]$$

discount factor

$$V(s) = \max_a \sum_{s' \in S} \{ R_{a,s}(s, s') + \gamma V(s') \}$$

$$\text{as } t \rightarrow \infty \quad \|V^{t+1}(s) - V^t(s)\|_2 \rightarrow 0$$

- algo. converges.

- Value iteration.

Monte-Carlo method

Q-learning

SARSA

Markov Decision Process (MDP)

\downarrow

Finite Horizon MDP

Infinite Horizon MDP

- Maximum expected utility.

- If we are in last, it doesn't matter what is the last state.

$$V^{T-1}(s) = \sum_{s'} \{ R_{a,s}(s, s') + \gamma \cdot V^T(s') \} \cdot P_{a,s}(s')$$

$\gamma < 1$

$\gamma = 0$, don't care about future rewards

$$V^T(s) = \max_a \sum_{s'} P_{a,s}(s') \{ R_{a,s}(s, s') + V^T(s') \}$$

this converges

→ we have only one state $|S|=1$, what policy should we use?
 → Maximum expected utility is used as optimal policy to use.

$$\pi^*(s) = \operatorname{argmax}_a R_a(s, s)$$

If we know that

$$R_a = 1 \text{ w.p. } M_a$$

$$= 0 \text{ otherwise}$$

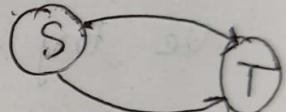
$$a^* = \operatorname{argmax} M_a$$

→ If we don't know M_a ,

we have 2 states S, T . There are some delays in the path. Which one would you choose?

- Multi-Armed Bandits (MAB).

- Check all possible routes, find the one which has minimum delay. This is a greedy approach.



Greedy Approach ():

for $i = 1 \rightarrow m k$

use action i } Exploration

observe a_i, R_i

update \hat{u}_i

for $i = m k + 1 \rightarrow T$

use action $i^* = \operatorname{argmax} \hat{u}_i$ } Exploitation

update \hat{u}_{i^*}

Issues: → A new smaller reward can dominate in first trial.

$$U_1 = 0.9, U_2 = 0.8$$

$\hat{U}_1 = 0 \rightarrow \hat{U}_2 = 1$

$$m = T/k$$

- using Round Robin

One every trial,

$$\text{loss} = \left[\frac{0.1 \times T}{2} \right]$$

Exploitation \rightarrow Et Greedy Algorithm

w.p. $(1 - \epsilon)$ use greedy strategy

w.p. ϵ select an action randomly

\rightarrow How to set ϵ ? learning rate η_t (w.p. $1/k$)

We try to find $\epsilon_t = O(1/t)$

as time increase, ϵ_t decreases.

\rightarrow How do we compare diff. algorithms?

Regret:

Expected loss in reward

$$\sum_{t=1}^T (U_a^* - U_a)$$

$\rightarrow O(T \log T) \cdot \Omega(\log T)$

\rightarrow No algorithm can do better than this.

Upper Confidence Bound (UCB) - Confident
also than prev

$\rightarrow O(\log T \cdot \sqrt{\log T})$

for
use (pull arm) action t
update \hat{m}_i^*

UCB

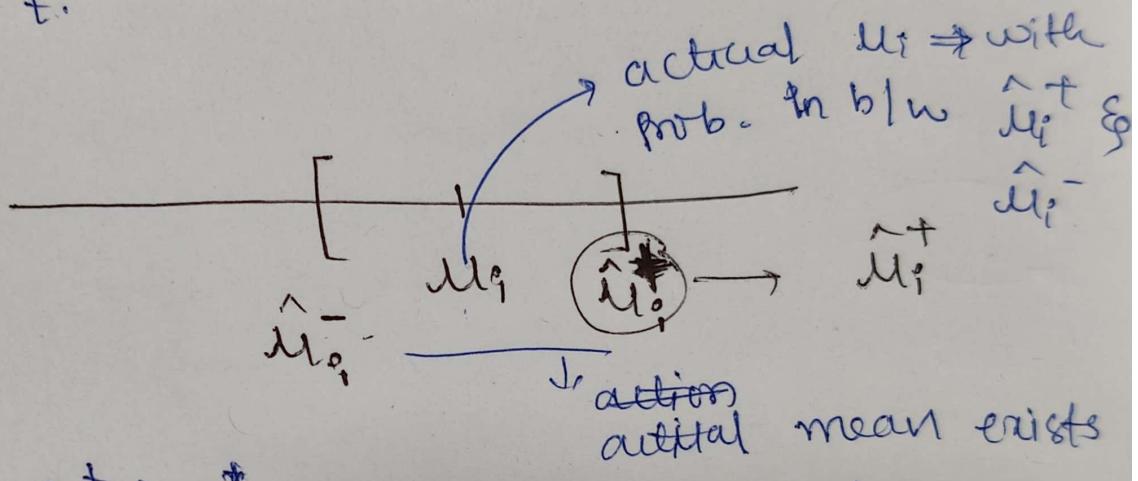
for $t = k+1 \rightarrow T$

$$\text{pull } a^* = \arg\max \hat{m}_i^t + \sqrt{\frac{2 \ln t}{N_{i,t}}}$$

empirical
mean

$\rightarrow N_{i,t}$

$\rightarrow N_{i,t}$ = Number of times
action/arm is used till
time t .



when $a^* \neq a^*$,
we will add reward.

- $O(\log T)$. \rightarrow since need to pull 'T' times.
no better also. to pull $< T$ times.
- $N_{i,t}$ is low, the interval will be high.
 \downarrow goes $(\log T)$ times.