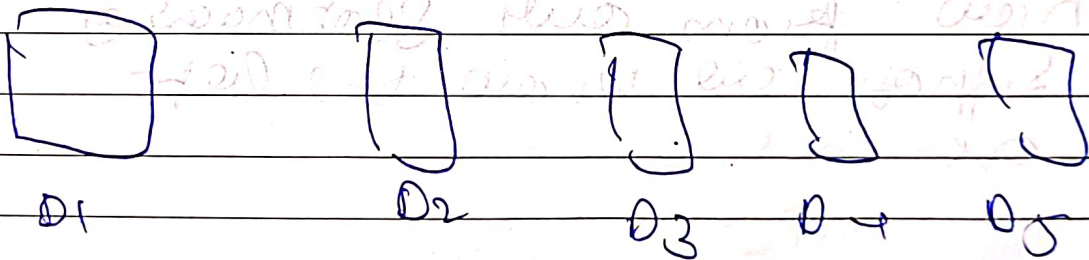


Reinforcement Learning

we need multi armed Bandit's problem

Here Bandit's is referred to a slot machine

So we have many different slot machines (eg 5)



each machine has a different distribution. and we need to figure out which has the best distribution

we don't know in advance which one is best

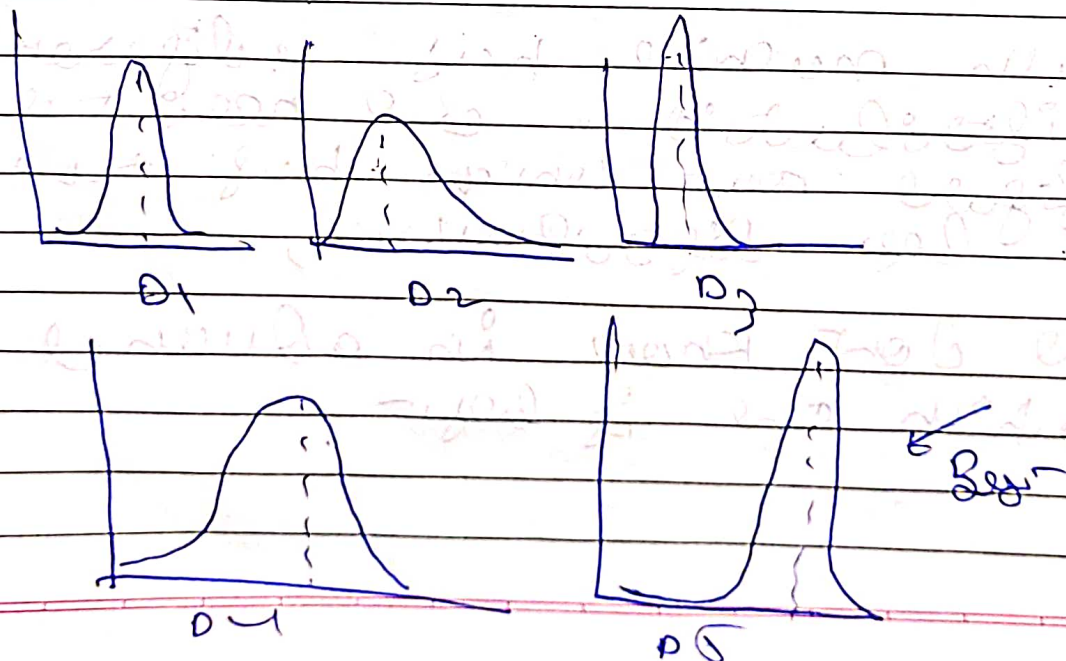
now every time we
use a machine we lose
money ~~all~~ so we need
to figure this out first

If we don't explore enough
then we will think a suboptimal
machine is optimal But we explore
for too long we lose money

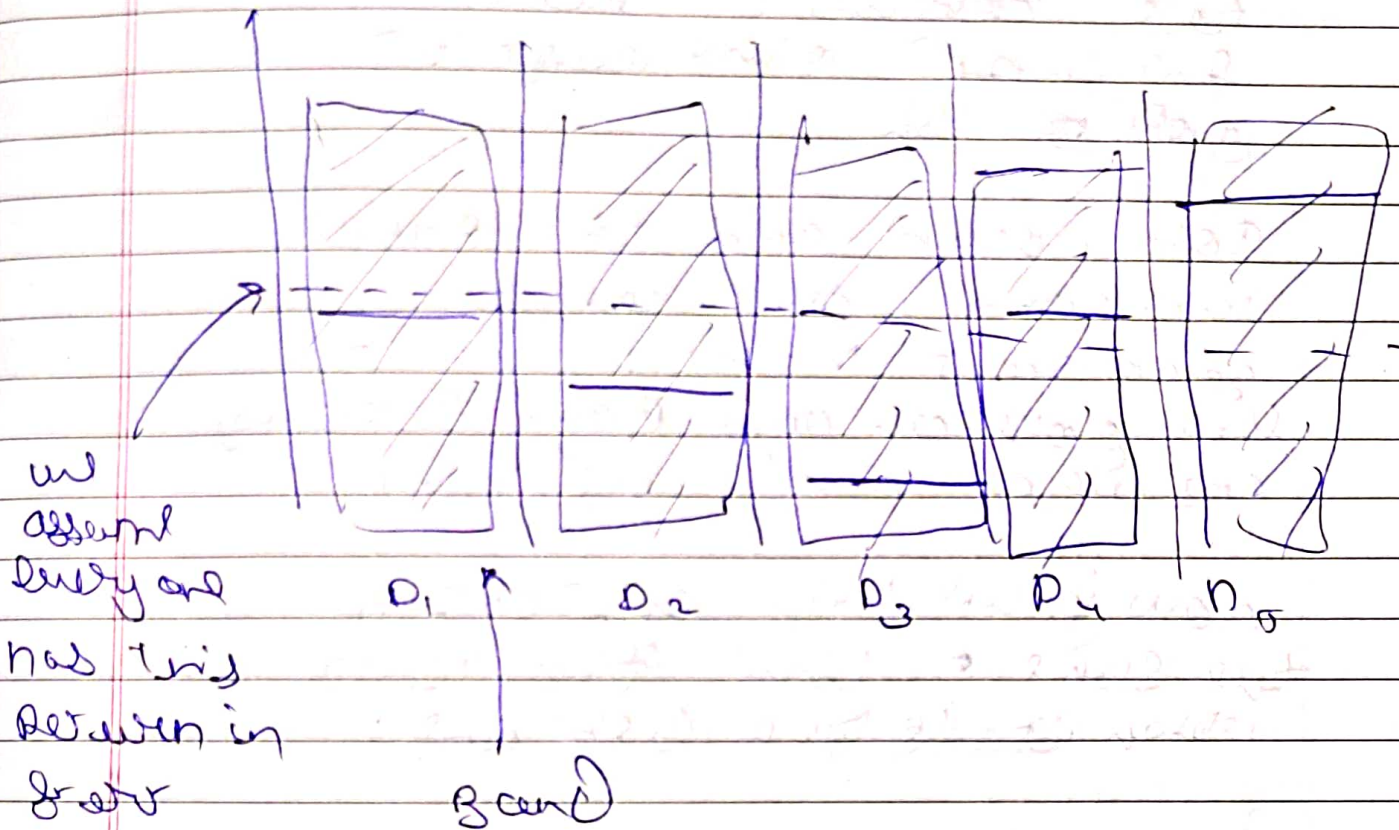
~~Time~~ on

→ upper confidence Bound

now from only slot machine
suppose we know the best
one



lets make a graph



we make the confidence Band such that the actual returns with a very high level of probability will be in the Band

Since both the machine with the higher confidence Band (as for all are same)

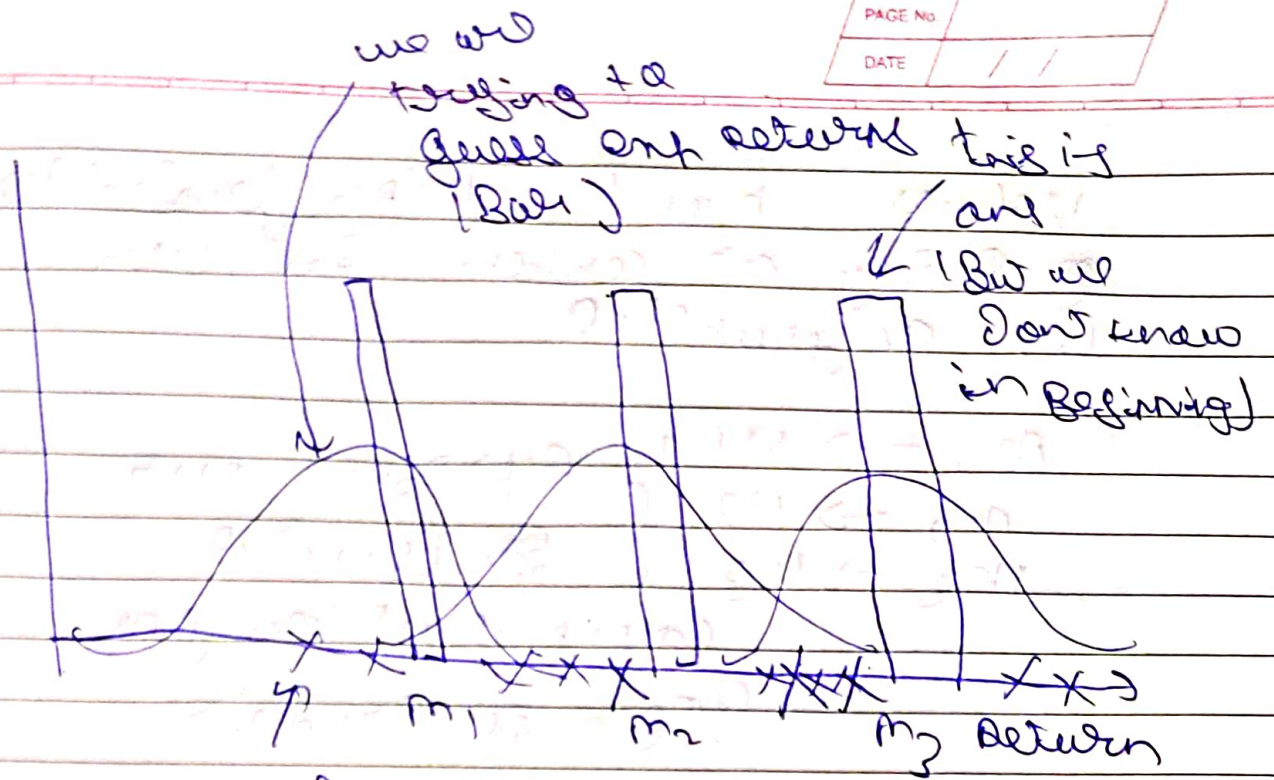
now if we test it and
it goes wrong (Below
the --- line) then
the ~~considered~~ ~~band~~ line
~~indicates otherwise it~~
~~goes up~~

goes down and the band
increases and if it
is correct then the ---
line goes up and band decreases
(as we are more confident)

now if we keep on doing
the same time after time
then it is the left end.

→ Implication Sampling

we will consider the
similar sample as
before.



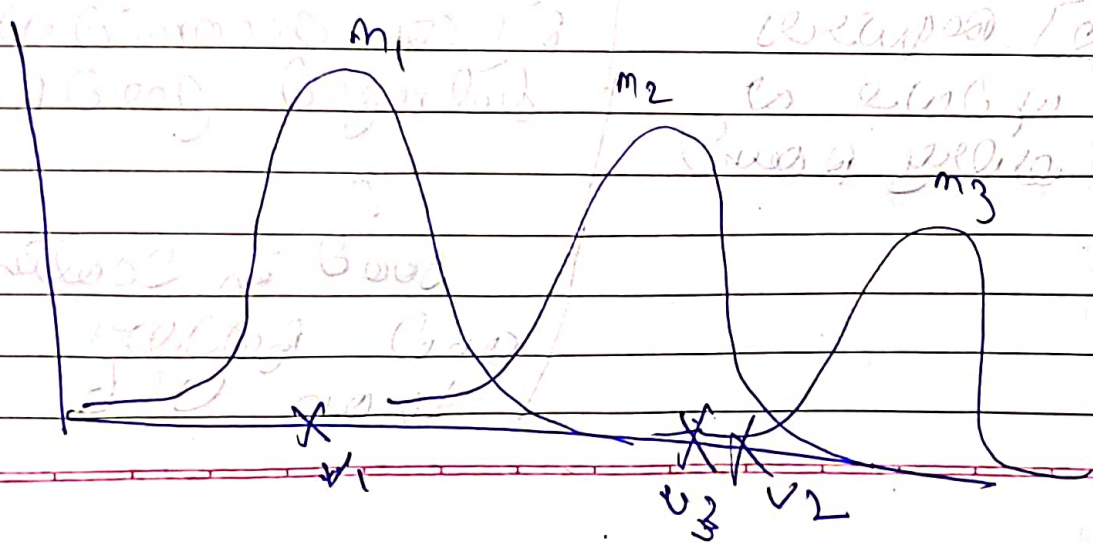
trial rounds

we will consider 3 Bandits

[we are not trying to guess distribution behind machine]

we are guessing output returned

now we will Pull some Random values from the distribution



values have high probability of being at center but can be anywhere

$m_1 \rightarrow v_1$
 $m_2 \rightarrow v_2$
 $m_3 \rightarrow v_3$

Represent our
 In the
 configuration for
 Logic of Bandit

So based on only values v_3 has max Return then we try m_3 then it's different solution sets and we full Repeat the process

UCB	Thompson Sampling
1) Deterministic	1) Probabilistic
2) requires update at every round	2) can accommodate delayed feedback ↗ good in scaling and better than UCB