Agenda:
- Bayesian Regret
- Markov Decision processes

---

Bayesian regret:
So far we studied "Frequentist" regret. For an environment $\nu$, and time step $n$, we studied $R(\nu, n)$.
we studied random regret, sometimes we took an expectation with respect to part of randomness in the reward. sometimes we took expectation with respect to all the randomness in the problem.

In the end, our regret bounds hold for all $\nu \in \mathcal{E}$.

---

Now consider the case that some one tells us $\nu \sim P_\nu$

and also tells us, given $\nu$, what is the exact distribution

of reward.

For example, assume for each arm $i$, $\mu_i \sim Beta(\alpha_i)$
and the reward is from $r_i \sim Bernoulli(\mu_i)$
Now we can ask, what is $\underset{\nu \sim P_\nu}{E}\left[R(\nu, \pi)\right]$ : Bayesian Regret

---

Posterior Sampling Reinfocement Learning (PSRL)

PSRL:
> First draw $\nu_1 \sim$ Prior,
> For $t = 1, \ldots$
>> Choose the optimal arm of $\nu_t$
>> Observe $X_t$
>> Update the posterior over $\nu_t$
>> Draw $\nu_{t+1} \sim$ Posterior

This simple algorithm usually gives us good Bayesian regret bound.

---

Example: Consider a 2-armed bandit with
where reward
of arm 1 is $N(\mu_1, 1)$
   arm 2 is $N(\mu_2, 1)$
where $\mu_1 \sim$ prior 1
       $\mu_2 \sim$ Prior 2



At time $t$ we compute the $\begin{cases} Posterior\ 1 \\ Posterior\ 2 \end{cases}$

draw $\mu_1 \sim$ Posterior 1 $\rightarrow \mu_1$ becomes 0.5
     $\mu_2 \sim$ Posterior 2 $\rightarrow \mu_2$ becomes $-0.2$
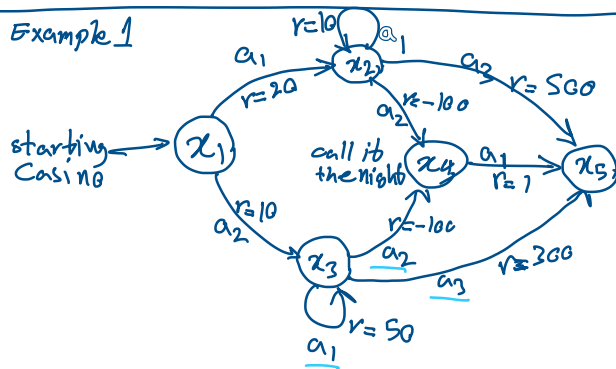Therefor arm 1 is optimal for this draw
we pull arm 1 and use the reward to update the
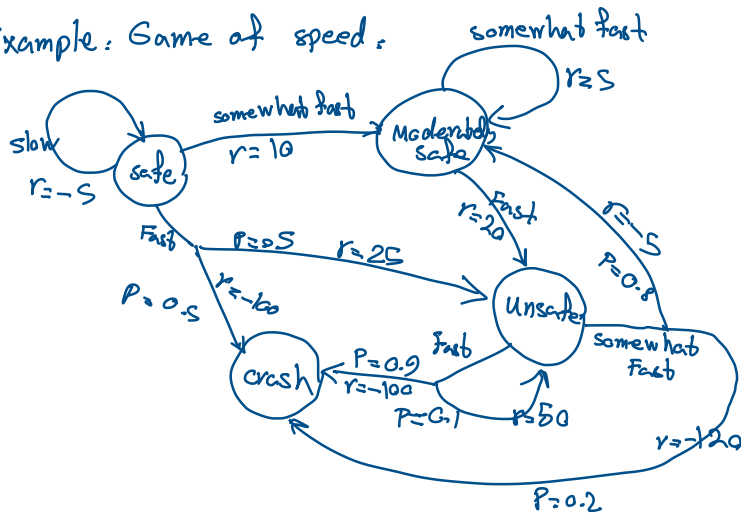Posterior of arm 1

# Markov Decision Process (MDP)

MDP is a controlled Markov process.

### Example 1



starting Casino $x_1$

$a_1$, $r=20$ → $x_2$
$r=10$ @ 1 (self loop on $x_2$)
$a_2$ $r=500$ → $x_5$
$a_2$ $r=-100$ → $x_4$
call it the night $x_4$
$a_1$ $r=1$ → $x_5$

$a_2$ $r=10$ → $x_3$
$x_3$ self loop $r=50$ $a_1$
$a_2$ $r=-100$ → $x_4$
$a_3$ $r=300$ → $x_5$

### Example: Game of speed.



slow $r=-5$ → safe
somewhat fast $r=10$ → Moderately safe
somewhat fast $r=5$ (self loop on Moderately safe)
Fast $P=0.5$ $r=25$ → Unsafe
Fast $P=0.5$ $r=-60$ → crash
Fast $r=20$ → Unsafe
$r=-5$ $P=0.8$
$P=0.9$ $r=-100$ → crash
Fast $P=0.1$ → crash
$P=50$
somewhat Fast $r=-120$
$P=0.2$

### Example: plain:     stochastic differential equation

$$\Rightarrow \quad \dot{x} = f(x, a) + d\beta$$

restless bandit where under each arm there is a Markov chain