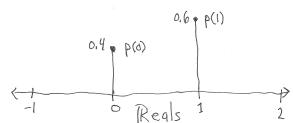


Random Variables, Expectations, Estimates, and a Learning Rule

Let $R \in \{0, 1\}$ $p(0) = 0.4$ $p(1) = 0.6$ $p(r) = \Pr[R=r]$



In general, consider a r.v. $R \in \mathbb{R} \subset \text{Reals}$ $\{R\} < \infty$

$R \sim P$ $P: \mathbb{R} \rightarrow [0, 1]$ s.t. $\sum_{r \in \mathbb{R}} p(r) = 1$

The expectation, or expected value of R is

$$\text{Defn: } E[R] = \sum_{r \in \mathbb{R}} p(r) \cdot r$$

Consider a sequence of r.v.s $R_t \sim P$ $t=1, 2, 3, \dots$

that are i.i.d. - identically, independently distributed

Consider the sequence of sample averages:

$$Q_{t+1} = \frac{R_1 + R_2 + \dots + R_t}{t} \approx E[R_t]$$

$$\lim_{t \rightarrow \infty} Q_t = E[R_t]$$

$$\text{Var}[Q_t] = E[(Q_t - E[R_t])^2] \propto \frac{1}{t}$$

Sample averages can be computed incrementally:

$$Q_{t+1} = Q_t + \frac{1}{t} [R_t - Q_t]$$

Which is a special case of our standard learning rule:

$$\text{NewEst.} = \text{OldEst} + \text{StepSize} \cdot [\text{Target} - \text{OldEst}]$$

$\underbrace{\phantom{\text{Target} - \text{OldEst}}}_{\text{error}}$

n-Armed Bandits

n actions $\in \{1, 2, 3, \dots, n\} = A$ $|A| = n$ "policy"

A sequence of r.v.s $A_t \in A$ $A_t \sim \pi_t$ \leftarrow not iid

$$t = 1, 2, 3, \dots \quad \pi_t : A \rightarrow [0, 1], \quad \sum_{a \in A} \pi_t(a) = 1$$

A sequence of r.v.s $R_t \in R \subset \text{Reals}$

$$p(r|a) = \Pr[R_t = r | A_t = a]$$

$$p : R \times A \rightarrow [0, 1]$$

$$\sum_{r \in R} p(r|a) = 1 \quad \forall a \in A$$

The value of action a :

$$q^*(a) = E[R_t | A_t = a] = \sum_{r \in R} p(r|a) \cdot r$$

Estimates

$$Q_t(a) \approx q^*(a)$$

greedy policy: $A_t = \arg \max_a Q_t(a)$

OR:

$$\pi_t(a) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} Q_t(a') \\ 0 & \text{otherwise} \end{cases}$$

ϵ -greedy policy:

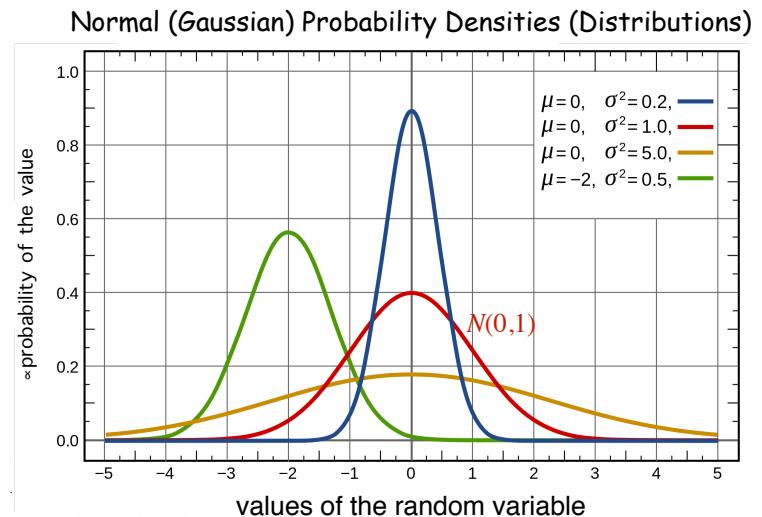
$$\pi_t(a) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{n} & \text{if } a = \arg \max_{a'} Q_t(a') \\ \frac{\epsilon}{n} & \text{otherwise} \end{cases}$$

Ten-Armed Testbed

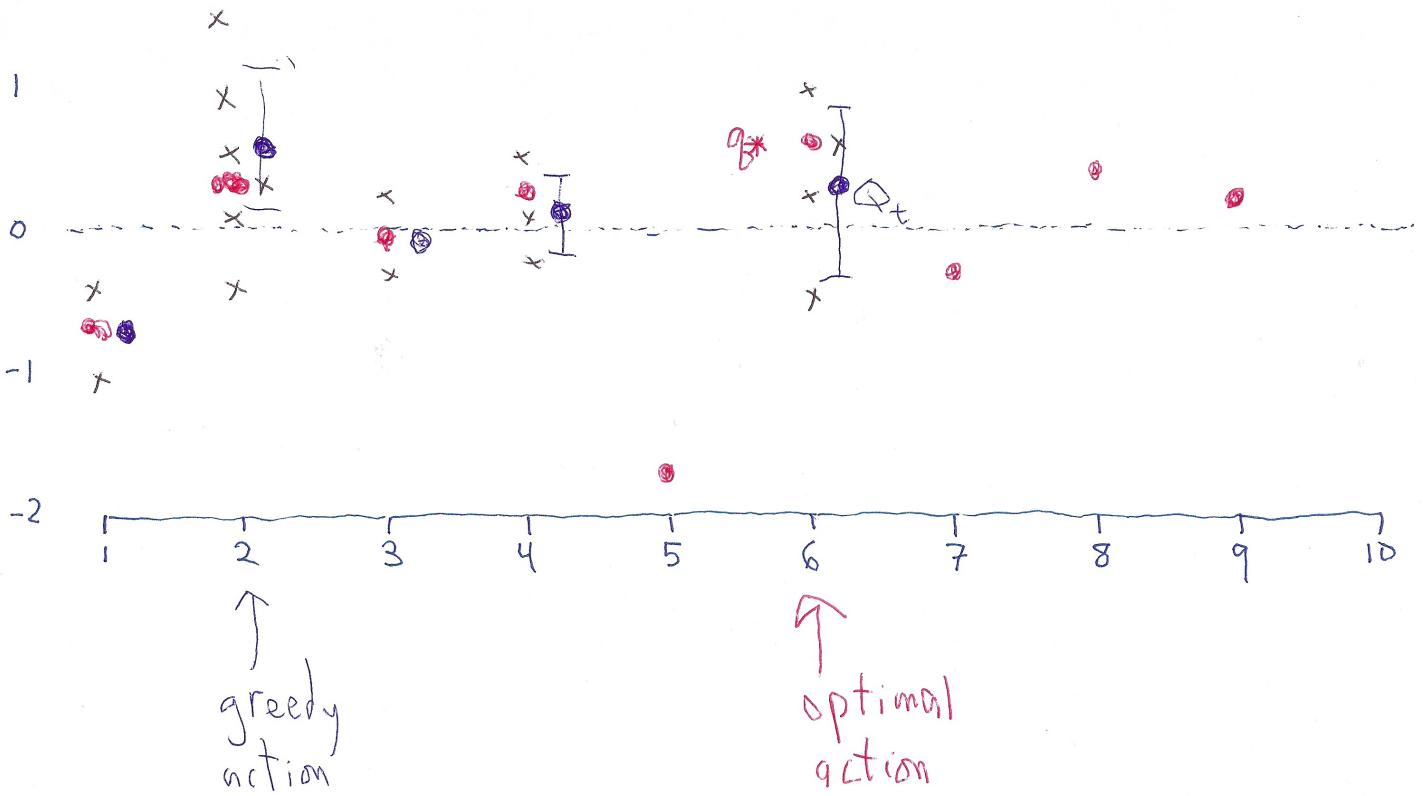
10 actions, $n=10$

$$q^*(a) \sim N(0, 1)$$

mean Variance



2



and

$$Q_{t+1}(A_t) = Q_t(A_t) + \alpha [R_t - Q_t(A_t)]$$

$$Q_{t+1}(a) = Q_t(a) \quad \forall a \neq A_t \quad a \in \mathcal{A}$$

What you have learned from bandits

- ϵ -greedy policies
- the difference between a sample, an estimate, and a true expected value - R_t, Q_t, q^*
- the difference between the greedy action and the optimal action
- a learning rule; how learning can be seen as computing an average in an incremental way
- seen a complete example of goal-seeking - both the problem and the solution methods
- seen a complete example of mathematical formalization of an AI problem & solution