$policy = solve(\underline{solver}, \underline{m})$     offline
                                                         calculation

$a = action(policy, s)$                                  online
                                                         calculation

# k-armed Bernoulli Bandits

$\theta \qquad p(w)$

$$p(\theta_i \mid o_1, o_2, \ldots o_{t-1}) = Beta$$

loop
    pull arm $i$ based on policy $\Leftarrow$
    observe $o_t$, get $r_t$
    $$P(\theta_i \mid o_1, o_2, \ldots o_{t-1}) = Beta(w+1, l+1)$$

---

$$\hat{Q}_i = \frac{w}{w+l} = \hat{\theta} \qquad \underset{i}{argmax}\ \hat{Q}_i$$
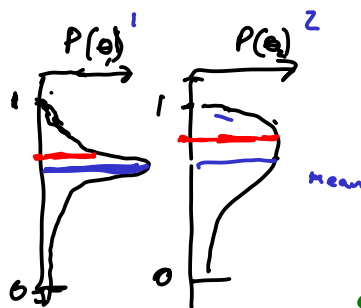
- $\varepsilon$-greedy:     wp. $1-\varepsilon$   choose $\underset{i}{argmax}\ \hat{Q}_i$
                 wp. $\varepsilon$   choose random

- softmax
$$P(i) = \frac{e^{\alpha \hat{Q}_i}}{\sum_i e^{\alpha \hat{Q}_i}}$$

- Interval

     · UCB

$$\underset{i}{argmax}\ \hat{Q}_i + c\sqrt{\frac{\ln \sum_i N_j}{N_i}}$$

- Thompson Sampling
     sample $\theta_i$ from $P(\theta_i \mid o_1, \ldots o_t)$
     choose $\underset{i}{argmax}\ \theta_i$

ad hoc

$P(\theta)^1$    $P(\theta)^2$    mean

$\hat{\theta} + 1$ standard dev

Bayesian   Expensive

(not expensive)

# Optimal Policies for Bandits — Finite Horizon h

## MDP

$$P(\theta \mid \ldots) = Beta(w+1, \ell+1)$$

$$S = \left\{ (w_1, \ell_1, \ldots, w_n, \ell_n) \right\}_{w, \ell \in 1..h}$$

exponential in h

$$A = \{ 1 \ldots k \}$$

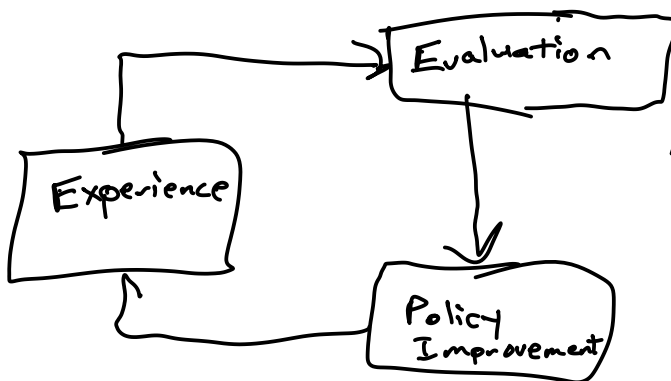$$R(s, i, s') = \begin{cases} +1 & \text{if } w_i \text{ was incremented} \\ 0 & \text{o.w.} \end{cases}$$

$T =$ increment win or loss according to belief probabilities

$$P(o_t = 1 \mid s) = \int P_{Beta(\ )}(\theta)\, \theta\, d\theta$$

## Gittins Index

~~~~~~~~~~~~~~~~~~~~

Exploration + Exploitation ← Bandits
Credit Assignment
Generalization

―――――――――――――――――――――

## RL                                    Offline or **Online**



often used in offline fashion with simulator

Model Based:  1. Learn MDP model from experience
              2. Solve MDP

Model Free:   Learn value or policy directly from experience

On-Policy

Off-Policy: Able to improve policy **without** new experience from that policy

# Model - Based RL

## Max Likelihood

$$N(s,a,s') \quad, \quad \rho(s,a)$$

$$N(s,a) = \sum_{s'} N(s,a,s')$$

$$T(s'|s,a) = \frac{N(s,a,s')}{N(s,a)}$$

$$R(s,a) = \frac{\rho(s,a)}{N(s,a)}$$

loop

    Choose a based on exp. strat.

    Observe $s', r$

    $N(s,a,s') \; ++$

    $\rho(s,a) \; += r$

    update $T, R$

    update $Q$ $\longleftarrow$ expensive . (Value Iteration)

    $s \leftarrow s'$