# Recap

| Method | Exploration/Exploitation | Optimal found? |
|---|---|---|
| A/B testing | 100% Exploration | YES — Through testing |
| Greedy | 100% Exploitation | NO — Local optimal |
| $\epsilon$-greedy | Explore at rate $\epsilon$ | YES — eventually |
| Softmax | Explore at softmax prob. rate | YES — better @ ranking than $\epsilon$-greedy |

**Note:** So far, greedy, $\epsilon$-greedy, and softmax look at the mean/proportion at each round <u>only</u>!
↳ There is no concern for <u>variability</u> or <u>uncertainty</u> of the reward.

• We can make optimization more efficient by accounting for the spread or variability of the reward distribution at each round.

<u>Two methods:</u>  1) Upper confidence bound (UCB)

2) Probability matching / Bayesian Sampling

[ Slide is wrong, correction: ]

$$UCB: \underset{k}{\arg\max} \left( r_{kt} + \sqrt{\frac{2 \log(t)}{n_{kt}}} \right)$$

So for UCB, instead of just updating $r_{kt}$ as in the previous methods, we now update $r_{kt} + \sqrt{\frac{2 \log(t)}{n_{kt}}}$ at each round.

# UCB in Practice

for $t = 1, \ldots, T$

choose arm with $\underset{K}{\arg\max} \left( r_{kt} + \sqrt{\dfrac{2 \log(t)}{n_{kt}}} \right)$

Return arm with highest value
↳ This is greedy in nature, and can be made more efficient by relying on the entire distribution of rewards.

# Randomized probability matching:

$\longrightarrow$ reward of arm k up to time t (estimated)

Value of interest: $P(r_{kt}$ is optimal$)$

$$P\left( r_{kt} = \underset{j \in 1, \ldots, K}{\max} \{ r_{jt} \} \right) = P_{kt}$$

we allocate units to arm $k$ with probability $P_{kt}$.

The probability we wrote up $(P_{kt})$ can be thought of in the following way:

$R = (r_1, r_2, \ldots, r_K) =$ True rewards for each arm

$\pi(R) =$ prior distribution on R

    a) If $r_j$ is a proportion, use a Beta distribution (or uniform $(0,1)$)

      Assume rewards are independent and each distributed as $U(0,1)$, then

      $\pi(R) = \overset{K}{\underset{j=1}{\prod}} 1 = 1$

    b) If $r_j$ is continuous, use a $N(\mu_j, 1)$ distribution where $\mu_j =$ the initialized reward for arm j.

      $\pi(R) = \overset{K}{\underset{j=1}{\prod}} \dfrac{1}{\sqrt{2\pi}} \exp \left\{ \dfrac{(r_j - \mu_j)^2}{2} \right\}$

      ( product of normals )

**Data:** The reward up to time $t$. That is $r_{1t}, ..., r_{kt}$

$$r_t = (r_{1t}, ..., r_{kt}) \leftarrow \text{estimated values up to round } t$$

*same as what is done in $\epsilon$-greedy.*

**Generating process:** $f(r_t | R)$

*rewards up to time $t$*     *true rewards*

Our value of interest can be thought of as the posterior distribution of $R$ given the rewards up to time $t$:

$(*)$    $P_{kt} = P\left( r_k = \max\{r_1, r_2, ..., r_k\} \mid (r_{1t}, ..., r_{kt}) \right)$

To get to $(*)$ is by calculating

$$P(R | r_{1t}, ..., r_{kt}) \propto \pi(R) f(r_{1t}, ..., r_{kt} | R)$$

Once we get $\nearrow$, we can then simulate say 1000 times and calculate the proportion of times $r_k$ was maximum as our value for $P_{kt}$.

1) rewards are proportions $\rightarrow$ use dirichlet-multinomial model
     (Dirichlet posterior)

2) rewards are continous $\rightarrow$ use normal-normal model
     (Normal posterior)

**For Proportions:**

**Algorithm:** Initialize and calculate rewards
   ↳ Normalize to probabilities for $t = 1, 2, ..., T$

Calculate Dirichlet posterior parameters for $P_{1t}, ..., P_{kt}$.

Simulate 1000 samples from Dirichlet and calculate

$$P_{kt} = P(r_k \text{ is max at time } t)$$

Pull arms with probabilities $P_{1t}, ..., P_{kt}$ and update prob.

**Algorithm:** Initialize and calculate rewards for $t = 1, 2, \ldots, T$

Calculate Normal posterior parameters for $r_{1t}, \ldots, r_{kt}$.

Simulate 1000 samples from Normal and calculate

$$P_{kt} = P(r_k \text{ is max at time } t)$$

Pull arms with probabilities $P_{1t}, \ldots, P_{kt}$ and update prob.

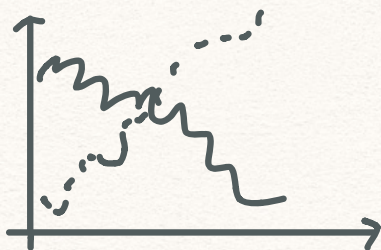$$\pi(R) = \prod_{j=1}^{k} f(r_j) \qquad x \sim U(0,1)$$
$$= \prod_{j=1}^{k} 1 \qquad f(x) = 1$$



# RPM (Randomized Prob. Matching)

- has been found to be the most efficient (in terms of iterations needed) to identify the optimal condition.
- This is viewed by looking at the optimality probability plots.



···· is optimal and converges to 1
—— converges to 0