

Last Time

Last Time

T? R?

- What is Reinforcement Learning?
- What are the main challenges in Reinforcement Learning?
 - Exploration & Exploitation
 - Credit Assignment
 - Generalization

Last Time

- What is Reinforcement Learning?
- What are the main challenges in Reinforcement Learning?
- How do we categorize RL approaches?

$$\pi_\theta \xrightarrow{\pi_\theta} \pi_{\theta'}$$

On Policy
Off Policy
Batch

Model Based TR
Model Free π_θ Q

Deep
Tabular

Last Time

Last Time

First RL Algorithm:

Last Time

First RL Algorithm:

Tabular Maximum Likelihood Model-Based Reinforcement Learning

Last Time

First RL Algorithm:

Tabular Maximum Likelihood Model-Based Reinforcement Learning

loop
 choose action a ↙
 gain experience
 estimate T, R
 solve MDP with T, R ↙

Guiding Questions

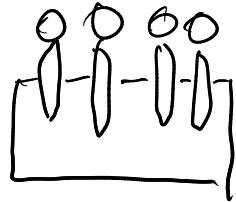
- What are the best ways to trade off Exploration and Exploitation

Bandits

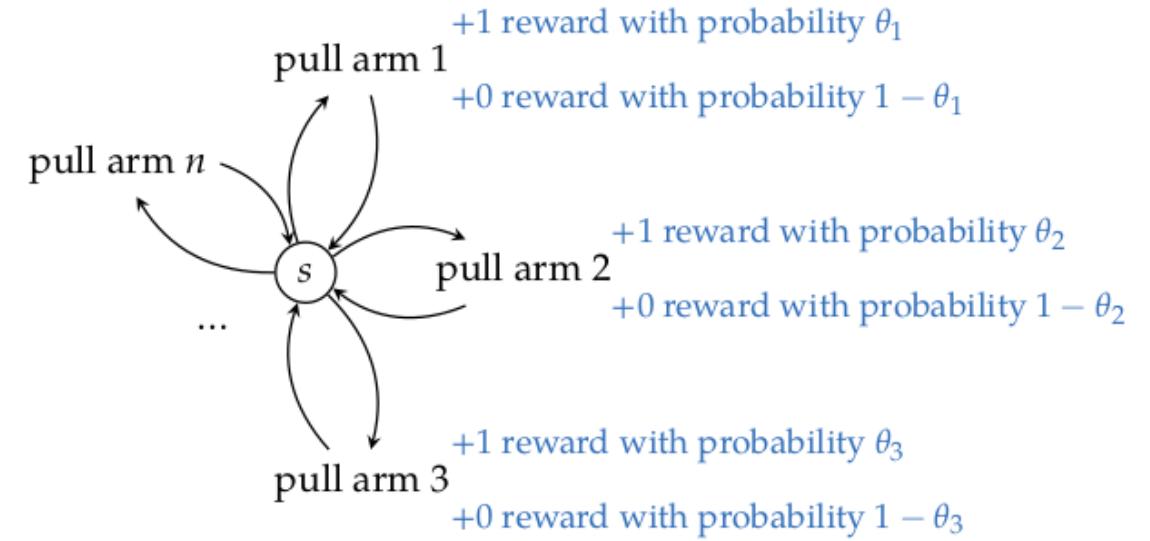


Bandits

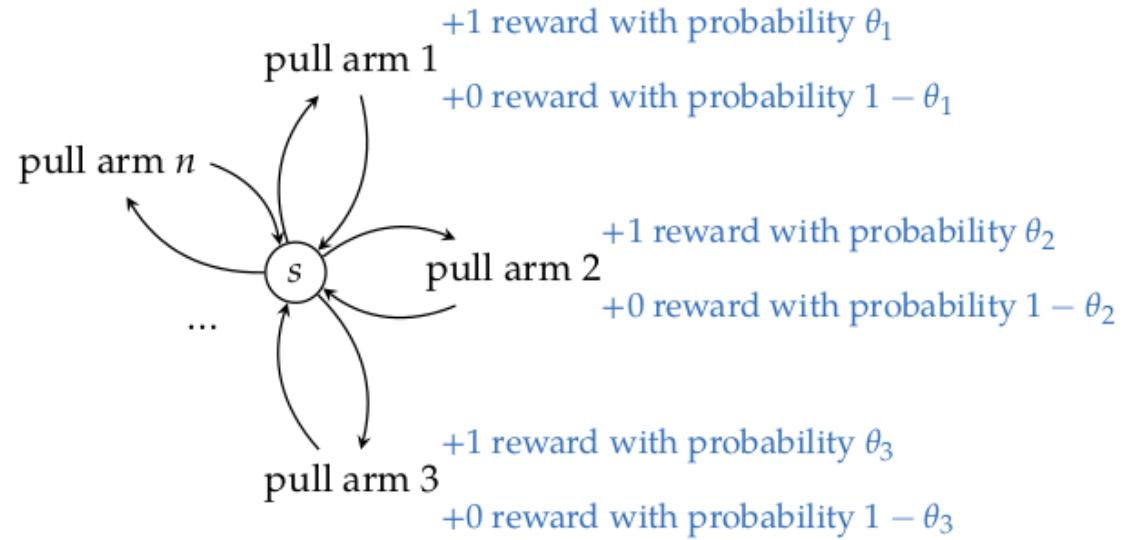




Bandits

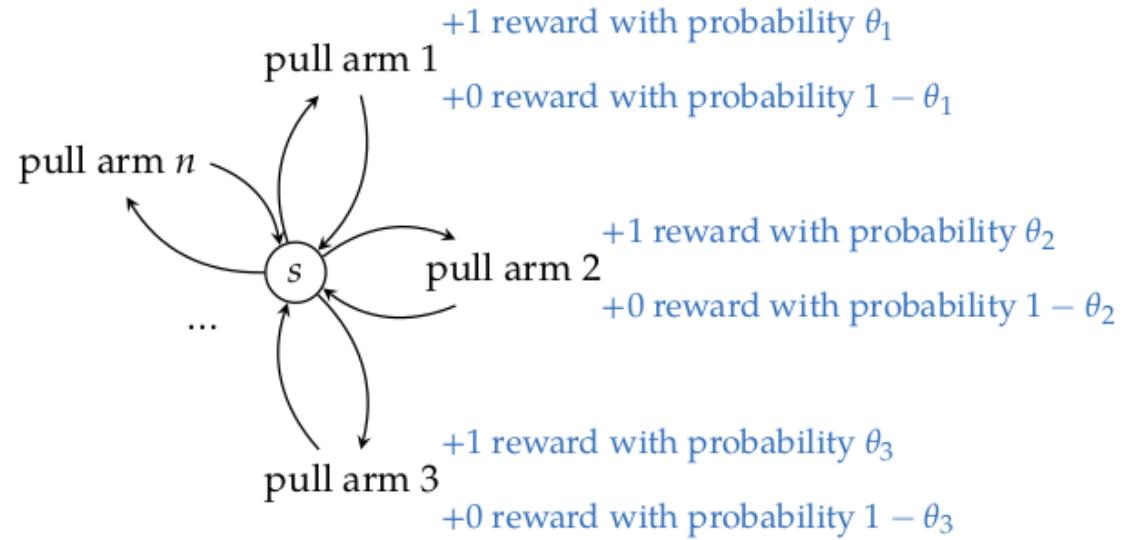


Bandits



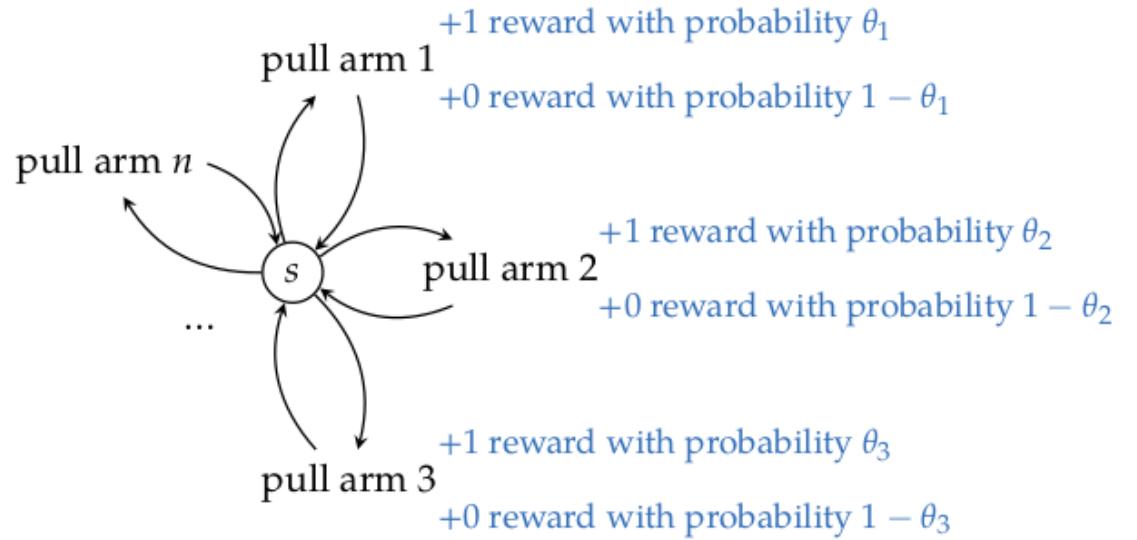
- Bernoulli Bandit with parameters θ

Bandits



- Bernoulli Bandit with parameters θ
- $\theta^* \equiv \max \theta$

Bandits



- Bernoulli Bandit with parameters θ
- $\theta^* \equiv \max \theta$

“According to Peter Whittle, “efforts to solve [bandit problems] so sapped the energies and minds of Allied analysts that the suggestion was made that the problem be dropped over Germany as the ultimate instrument of intellectual sabotage.”

Greedy Strategy

$$\rho_a = \frac{\text{number of wins}}{\text{number of tries}}$$

Choose $\operatorname{argmax}_a \rho_a$

Undirected Strategies

Undirected Strategies

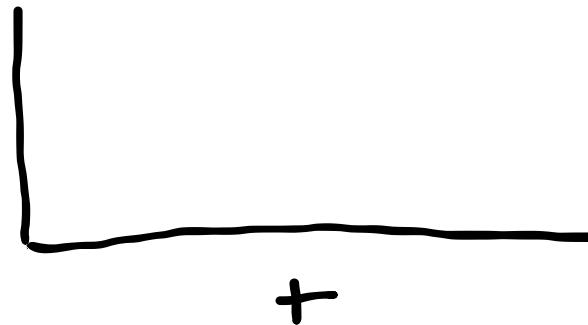
- Explore then Commit

Choose a randomly for k steps

Then choose $\underset{a}{\operatorname{argmax}} \rho_a$

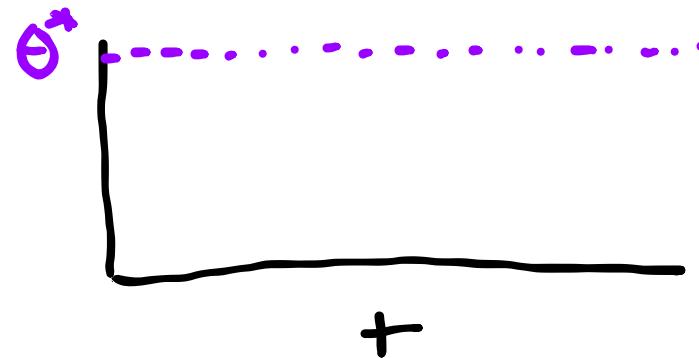
Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$



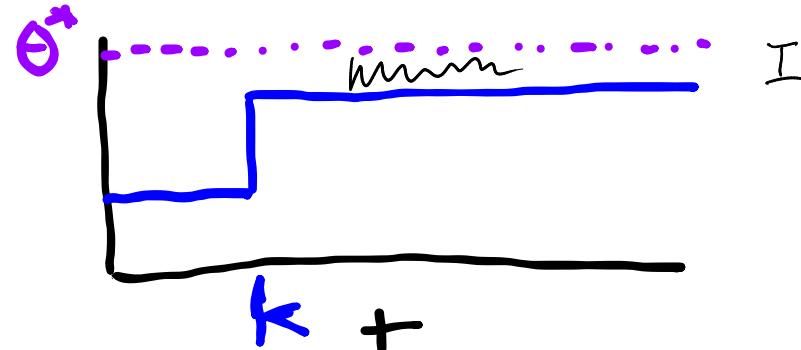
Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$



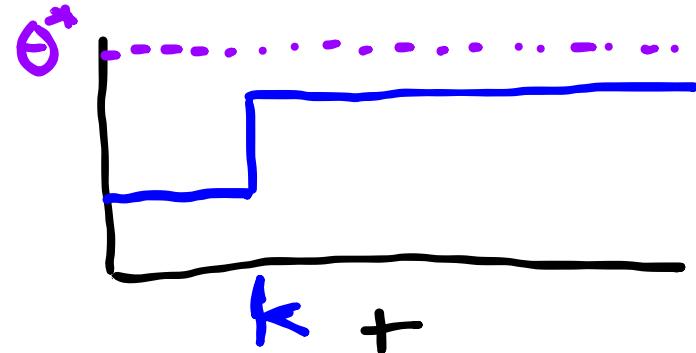
Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$



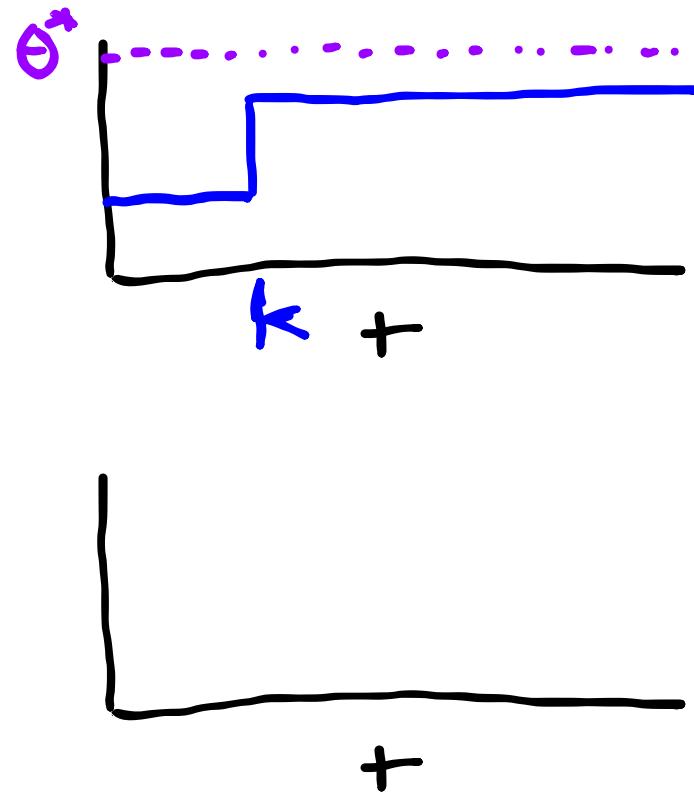
Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$
- ϵ - greedy
With probability ϵ , choose randomly
Otherwise choose $\operatorname{argmax}_a \rho_a$



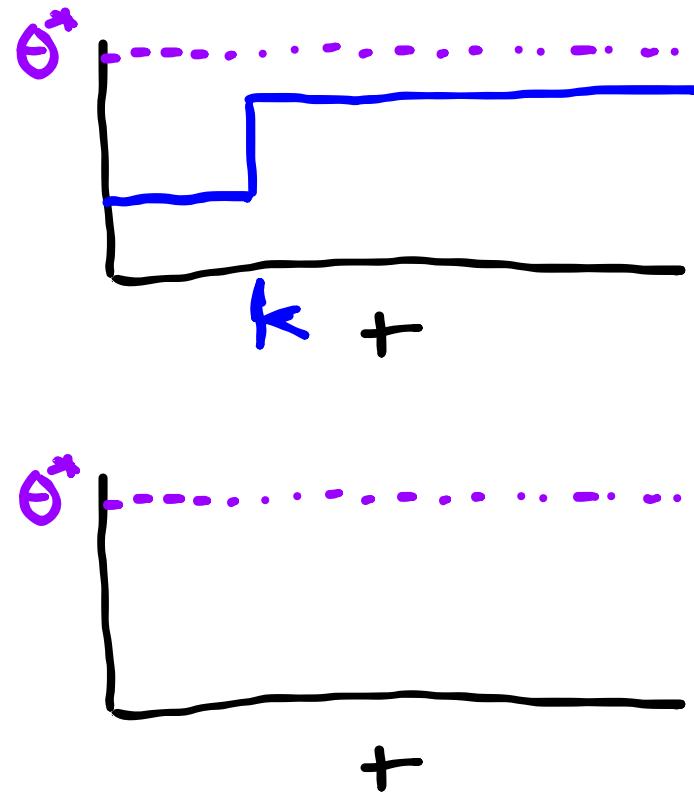
Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$
- ϵ - greedy
With probability ϵ , choose randomly
Otherwise choose $\operatorname{argmax}_a \rho_a$



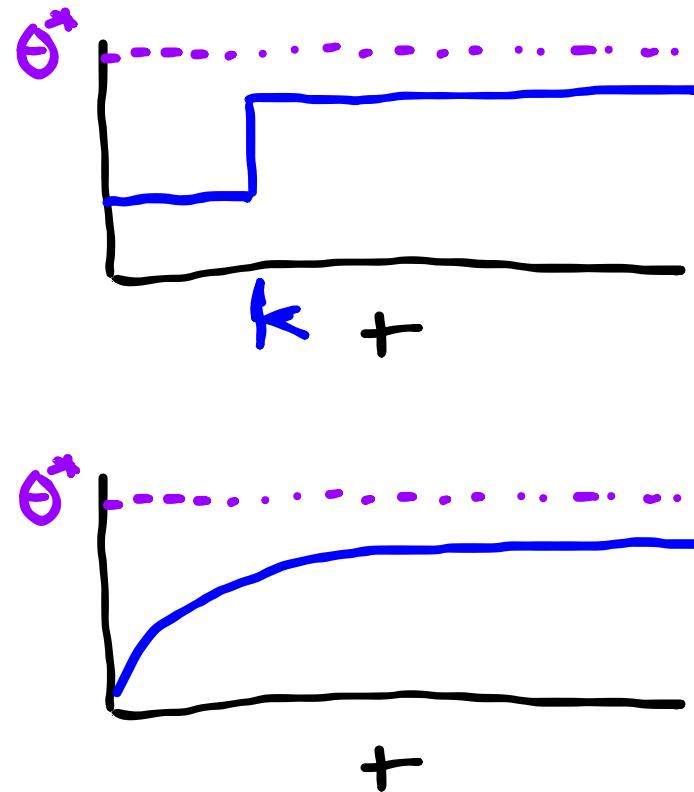
Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$
- ϵ - greedy
With probability ϵ , choose randomly
Otherwise choose $\operatorname{argmax}_a \rho_a$



Undirected Strategies

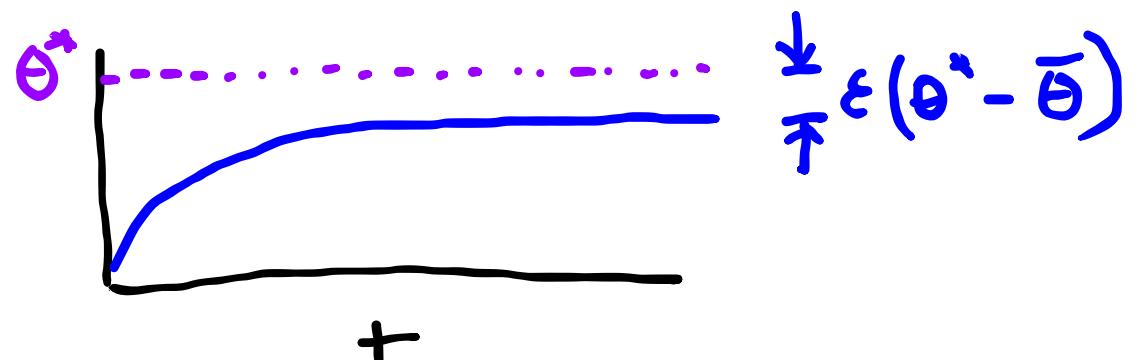
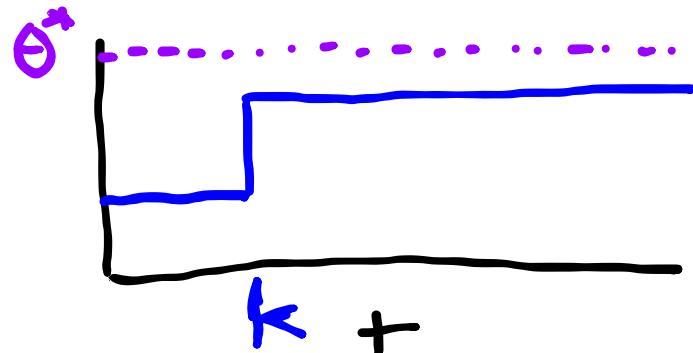
- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$
- ϵ - greedy
With probability ϵ , choose randomly
Otherwise choose $\operatorname{argmax}_a \rho_a$



Undirected Strategies

- Explore then Commit
Choose a randomly for k steps
Then choose $\operatorname{argmax}_a \rho_a$
- ϵ - greedy
With probability ϵ , choose randomly
Otherwise choose $\operatorname{argmax}_a \rho_a$

$$\begin{aligned}\epsilon &\leftarrow \alpha \epsilon \\ \alpha &\in (0, 1) \\ &\text{close to 1}\end{aligned}$$



Directed Strategies

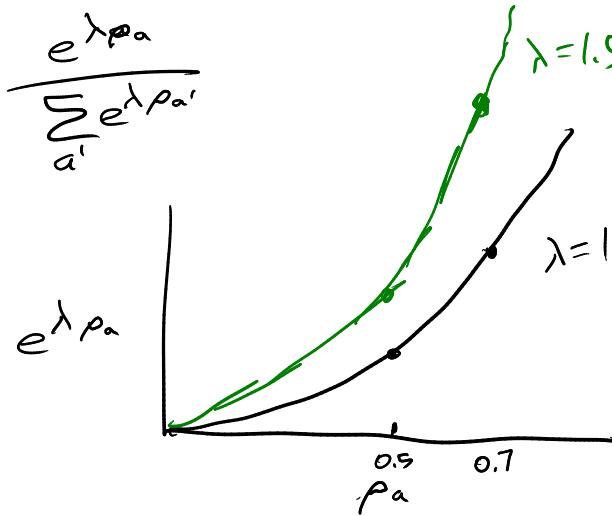


Directed Strategies

$$p(a) = \frac{e^{\lambda \rho_a}}{\sum_{a'} e^{\lambda \rho_{a'}}$$

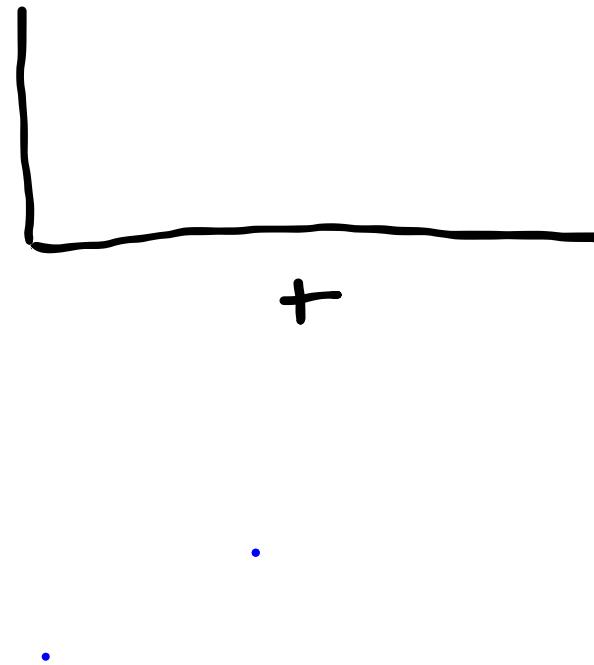
as $\lambda \rightarrow \infty$
looks like
 $\underset{a}{\operatorname{argmax}} p$

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$



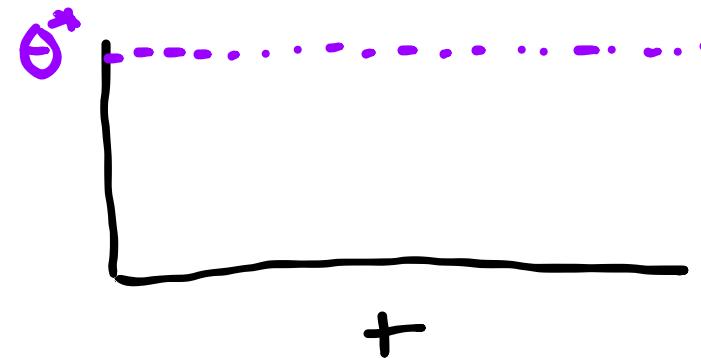
Directed Strategies

- Softmax
Choose a with probability
proportional to $e^{\lambda \rho_a}$



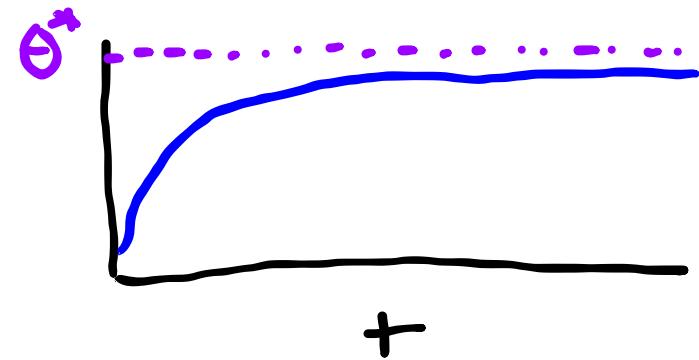
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$



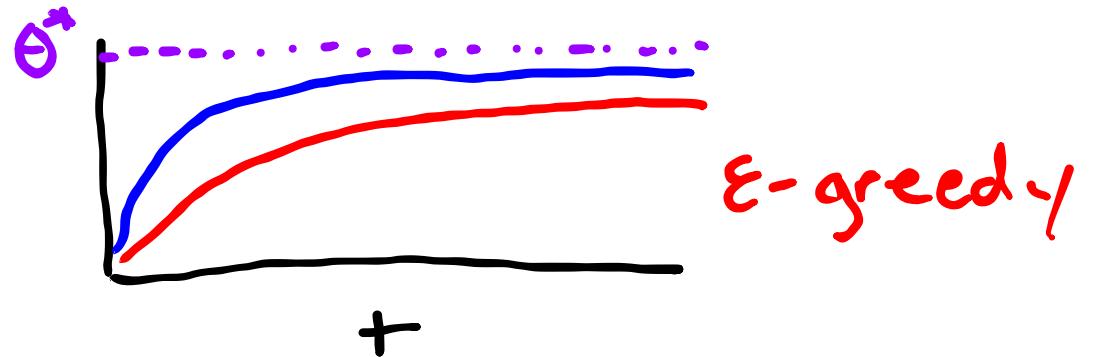
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$



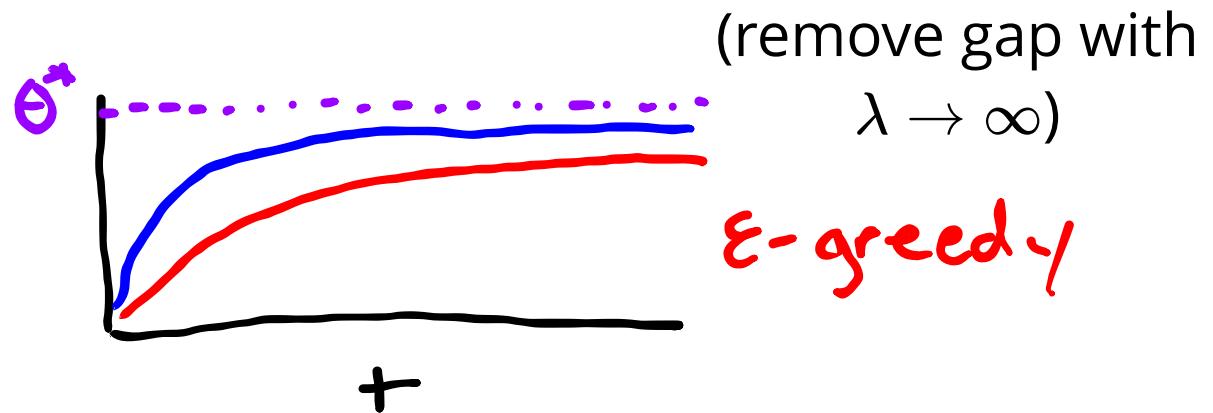
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$



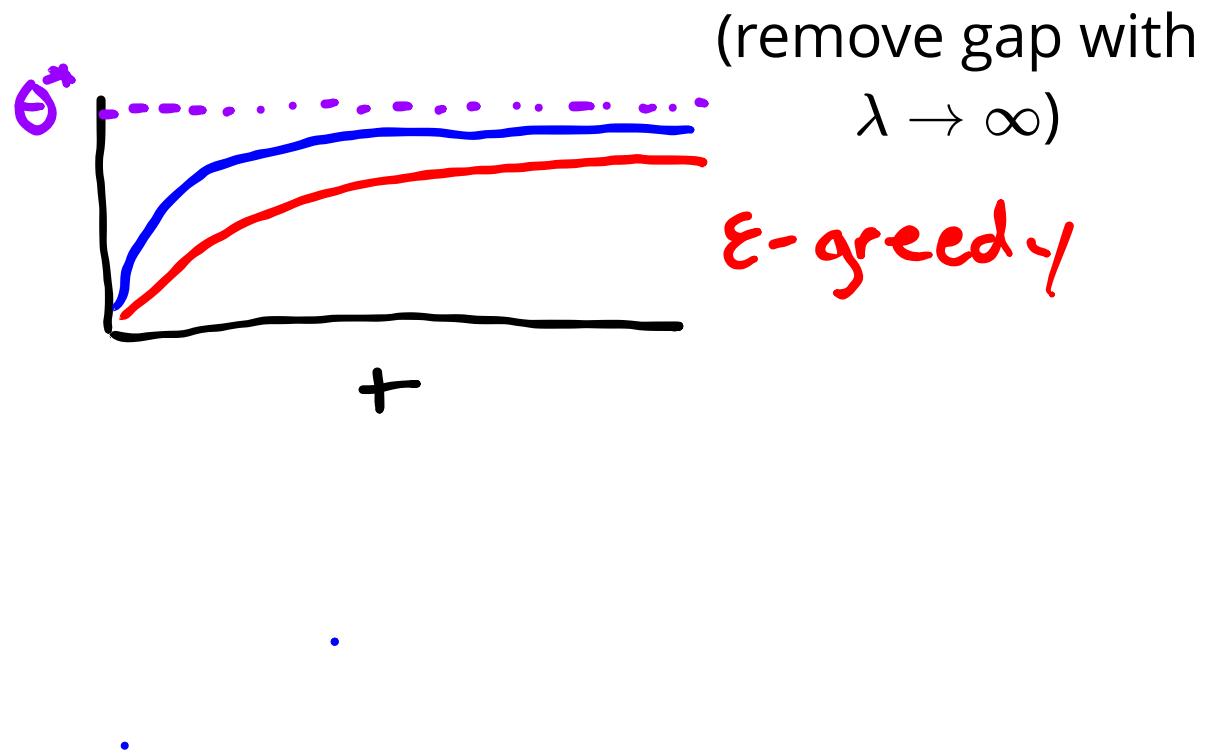
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$



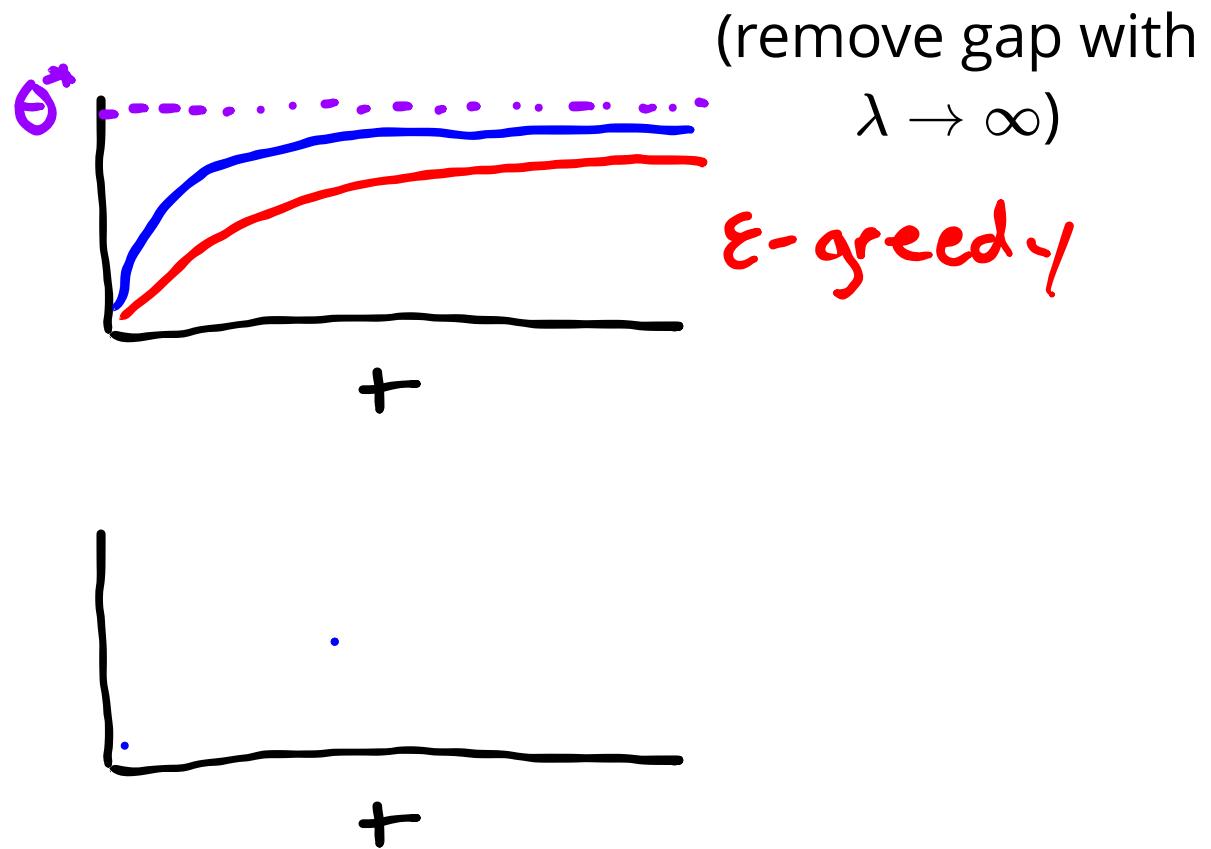
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$
- Upper Confidence Bound (UCB)
Choose $\underset{a}{\operatorname{argmax}} \rho_a + c \sqrt{\frac{\log N}{N(a)}}$



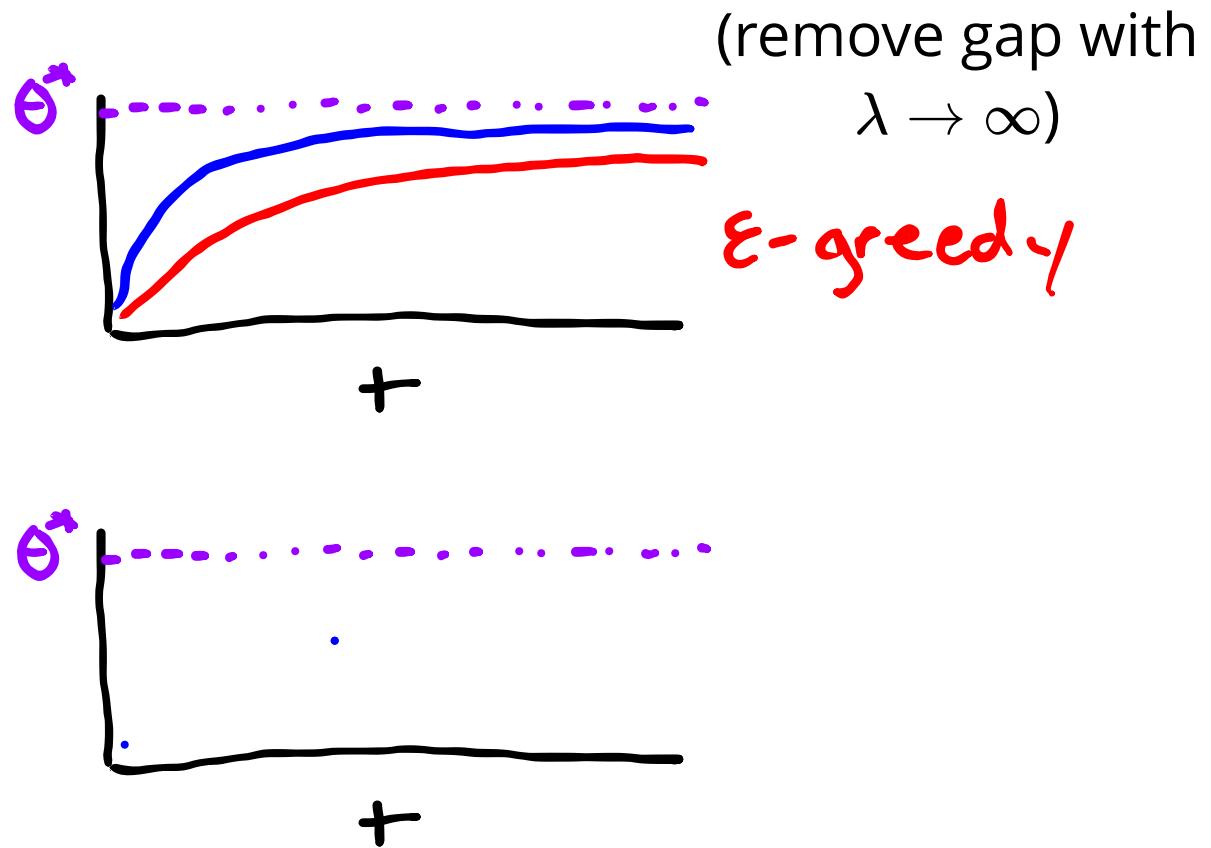
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$
- Upper Confidence Bound (UCB)
Choose $\underset{a}{\operatorname{argmax}} \rho_a + c \sqrt{\frac{\log N}{N(a)}}$



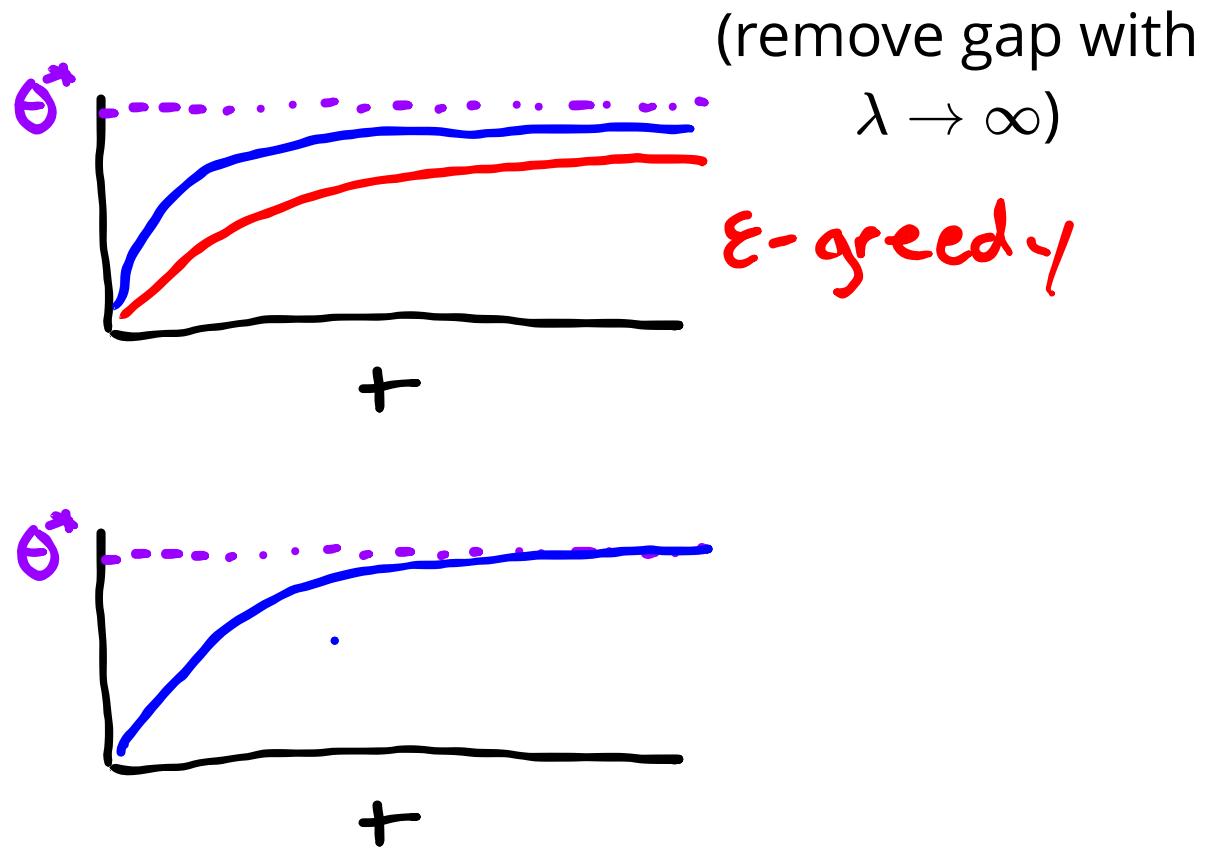
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$
- Upper Confidence Bound (UCB)
Choose $\underset{a}{\operatorname{argmax}} \rho_a + c \sqrt{\frac{\log N}{N(a)}}$



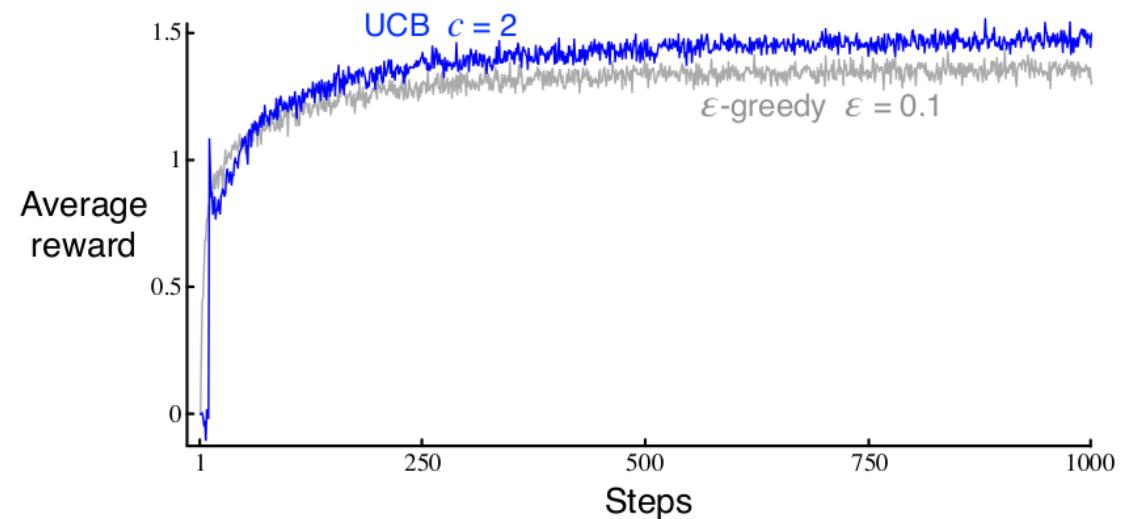
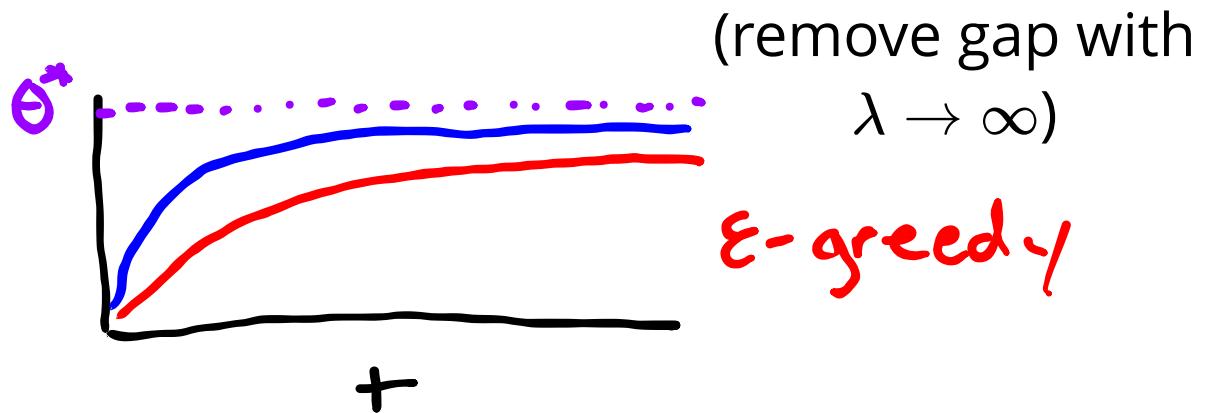
Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$
- Upper Confidence Bound (UCB)
Choose $\underset{a}{\operatorname{argmax}} \rho_a + c \sqrt{\frac{\log N}{N(a)}}$



Directed Strategies

- Softmax
Choose a with probability proportional to $e^{\lambda \rho_a}$
- Upper Confidence Bound (UCB)
Choose $\underset{a}{\operatorname{argmax}} \rho_a + c \sqrt{\frac{\log N}{N(a)}}$



Bayesian Estimation

Bayesian Estimation

Bernoulli Distribution

Bayesian Estimation

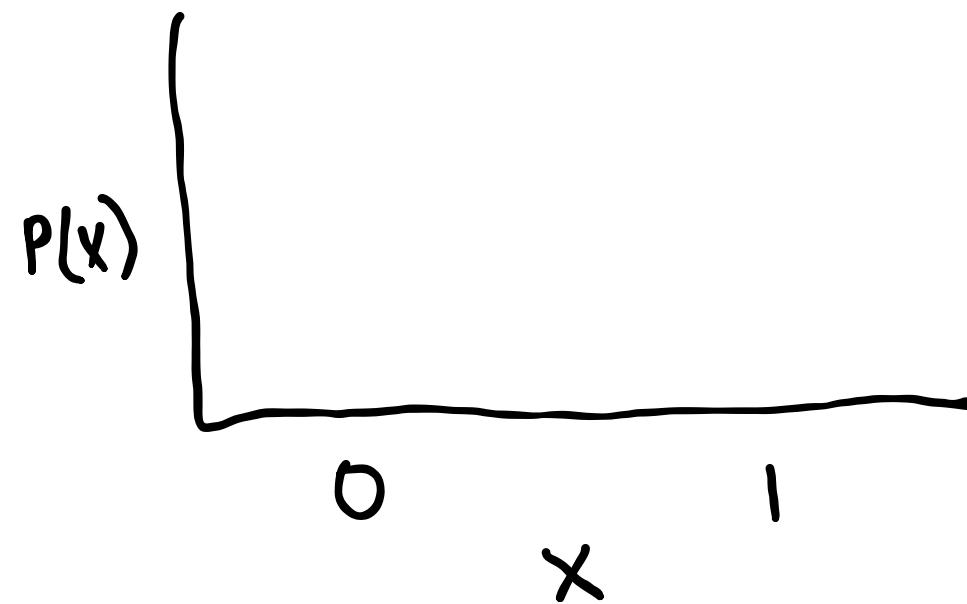
Bernoulli Distribution

$$\text{Bernoulli}(\theta)$$

Bayesian Estimation

Bernoulli Distribution

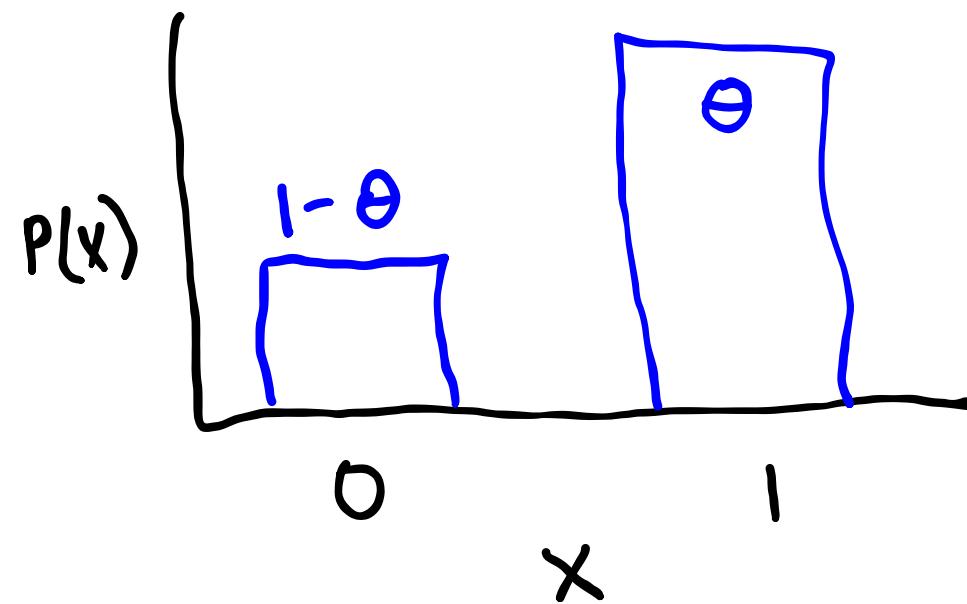
$$X \sim \text{Bernoulli}(\theta)$$



Bayesian Estimation

Bernoulli Distribution

$$\text{Bernoulli}(\theta)$$

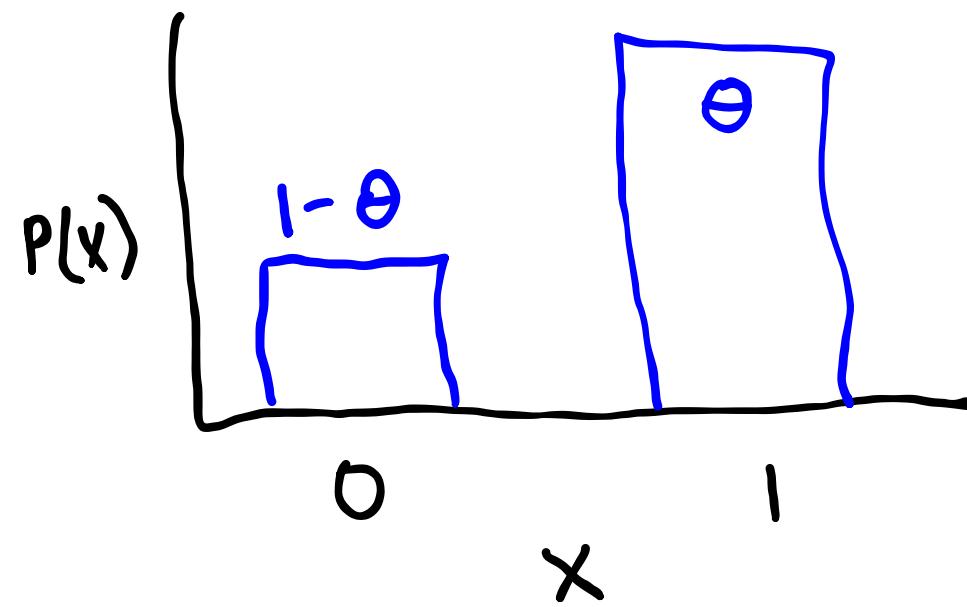


Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

Discussion: Given that I have received w wins and l losses, what should my belief (probability distribution) about θ look like?

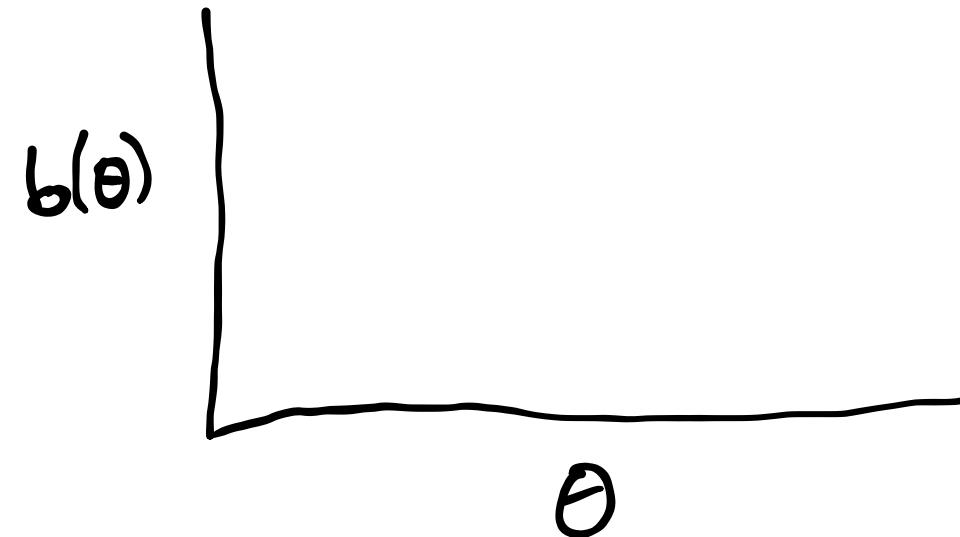
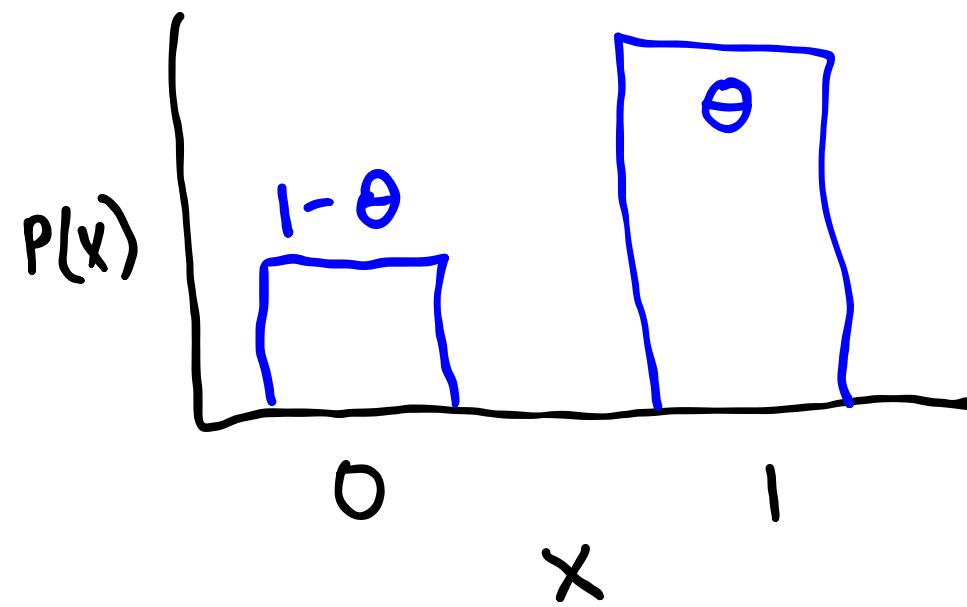


Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

Discussion: Given that I have received w wins and l losses, what should my belief (probability distribution) about θ look like?

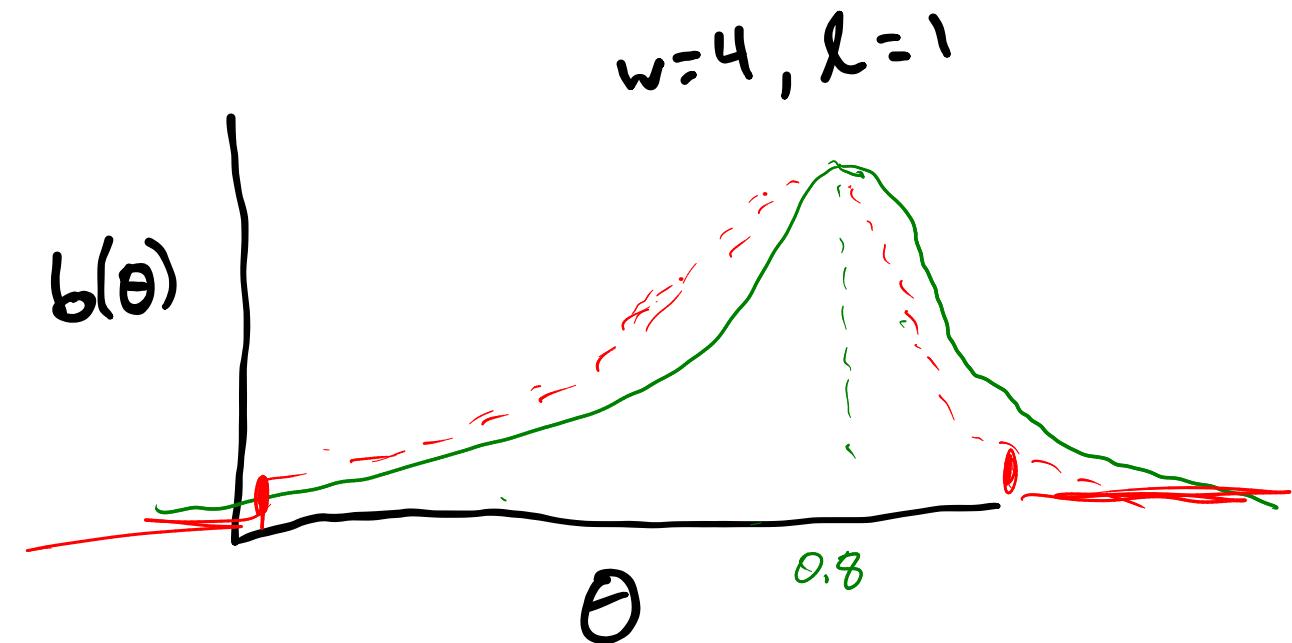
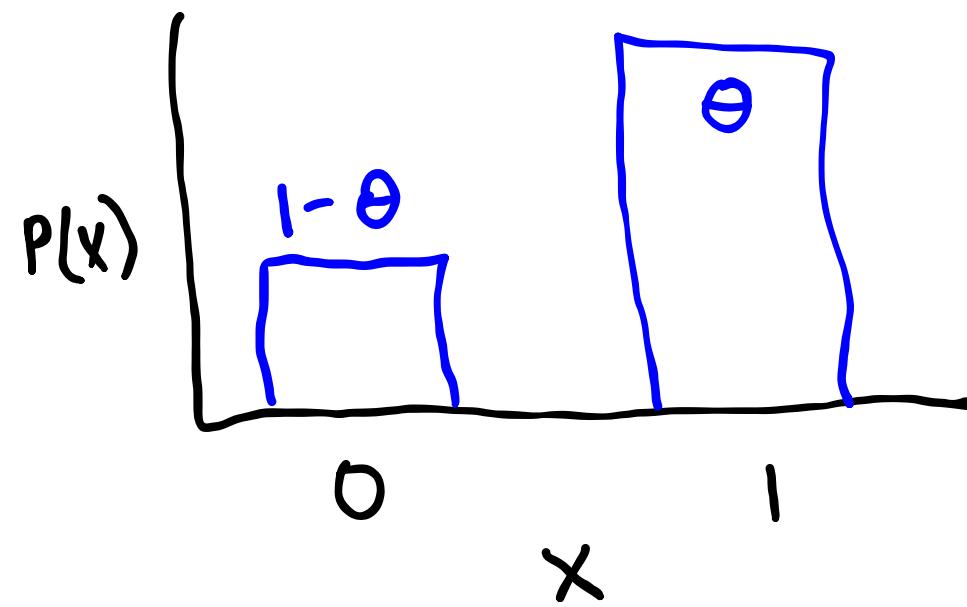


Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

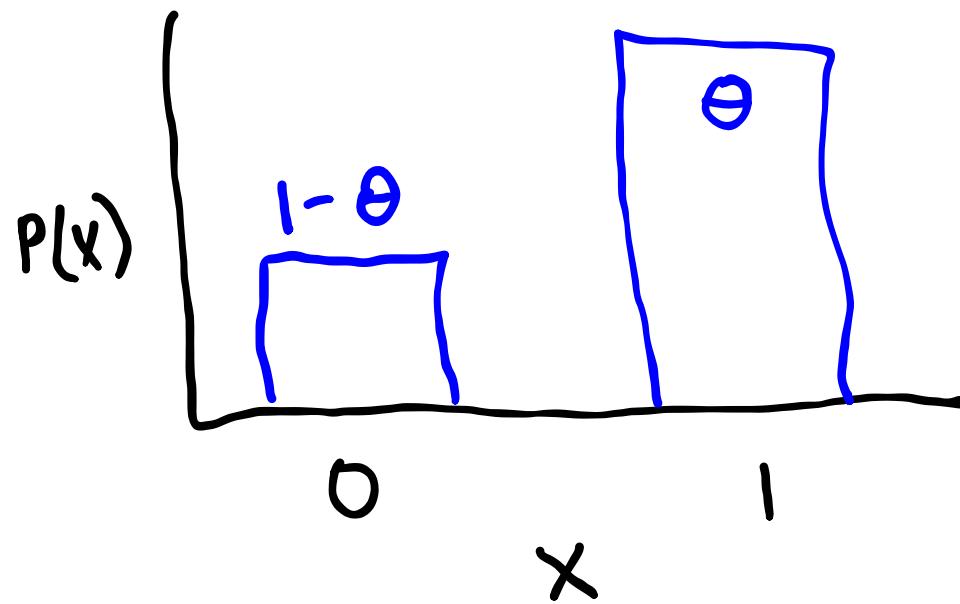
Discussion: Given that I have received w wins and l losses, what should my belief (probability distribution) about θ look like?



Bayesian Estimation

Bernoulli Distribution

$$\text{Bernoulli}(\theta)$$



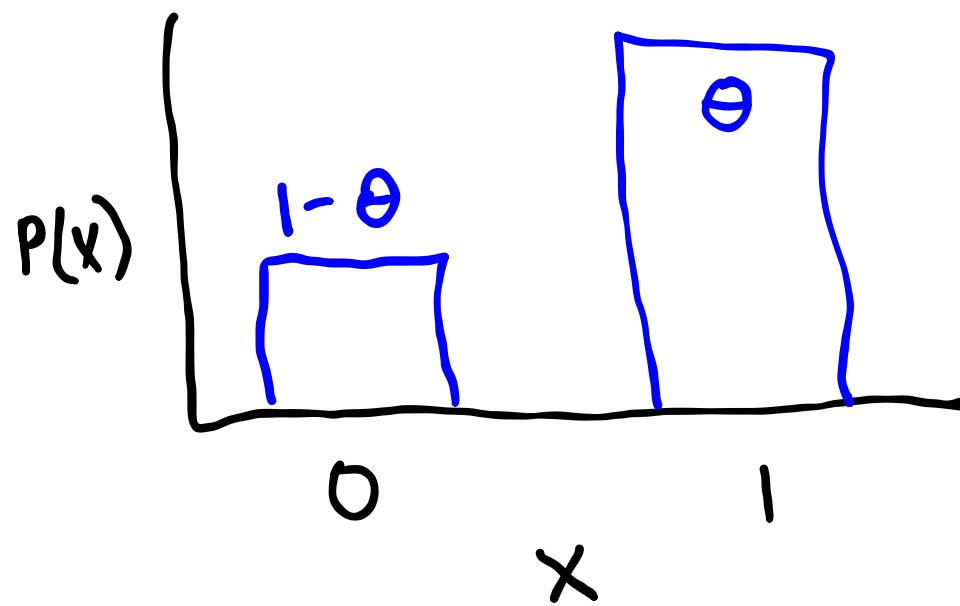
Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

Beta Distribution

(distribution over Bernoulli distributions)



Bayesian Estimation

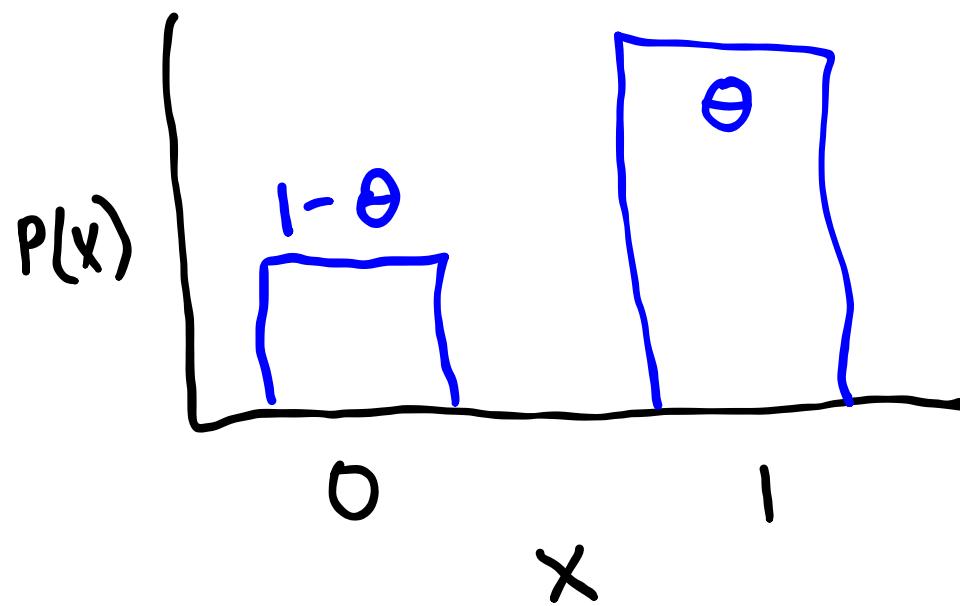
Bernoulli Distribution

$$\text{Bernoulli}(\theta)$$

Beta Distribution

(distribution over Bernoulli distributions)

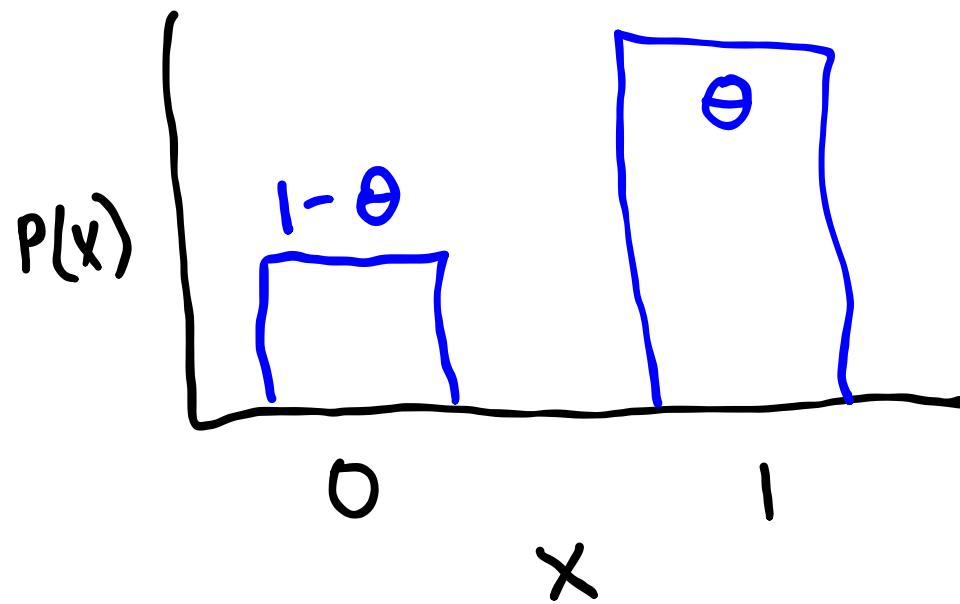
$$\text{Beta}(\alpha, \beta)$$



Bayesian Estimation

Bernoulli Distribution

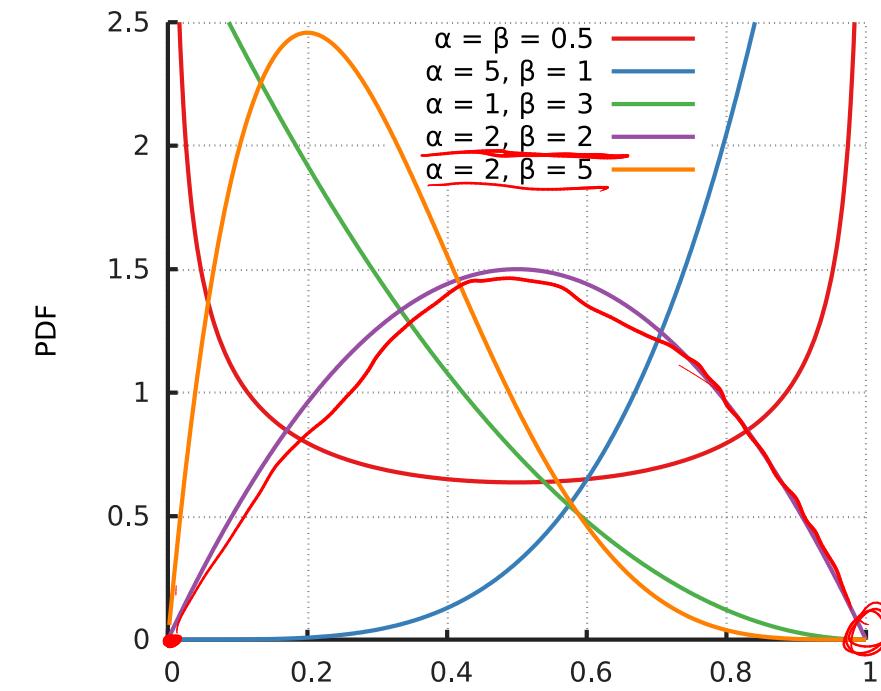
$\text{Bernoulli}(\theta)$



Beta Distribution

(distribution over Bernoulli distributions)

$\text{Beta}(\alpha, \beta)$



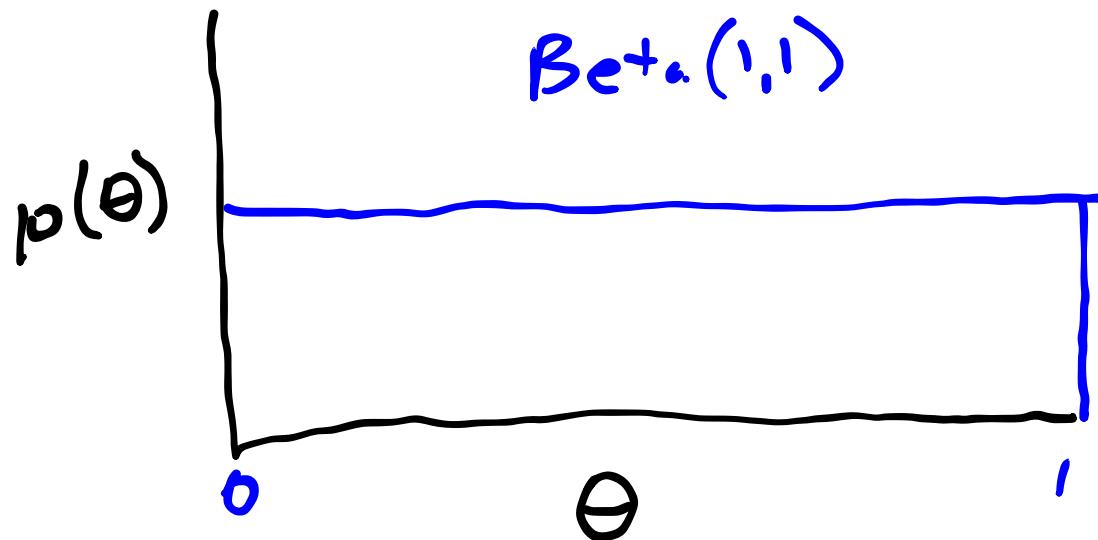
Bayesian Estimation

Bayesian Estimation

Given a $\text{Beta}(1, 1)$ prior distribution

Bayesian Estimation

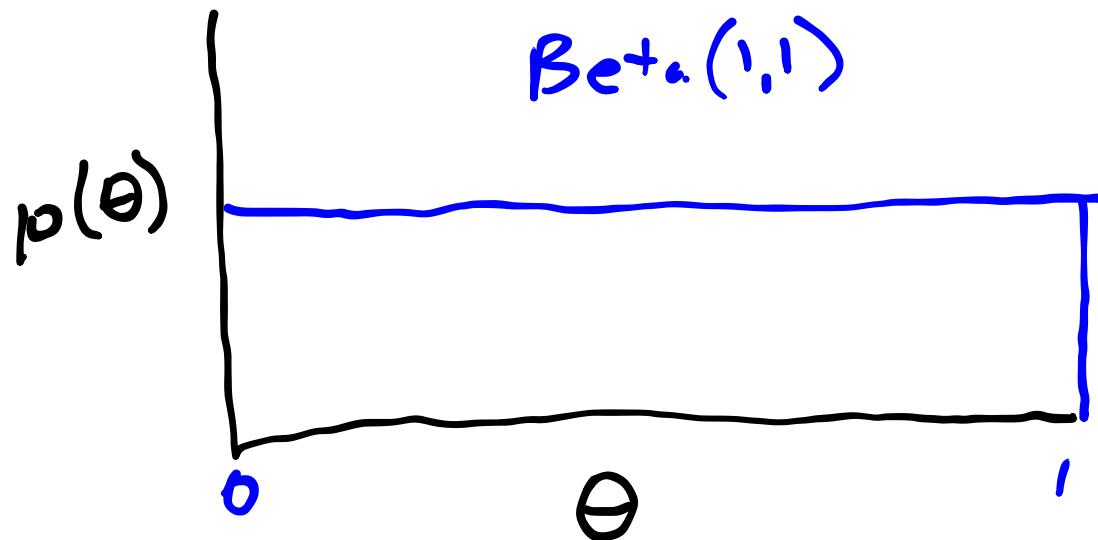
Given a $\text{Beta}(1, 1)$ prior distribution



Bayesian Estimation

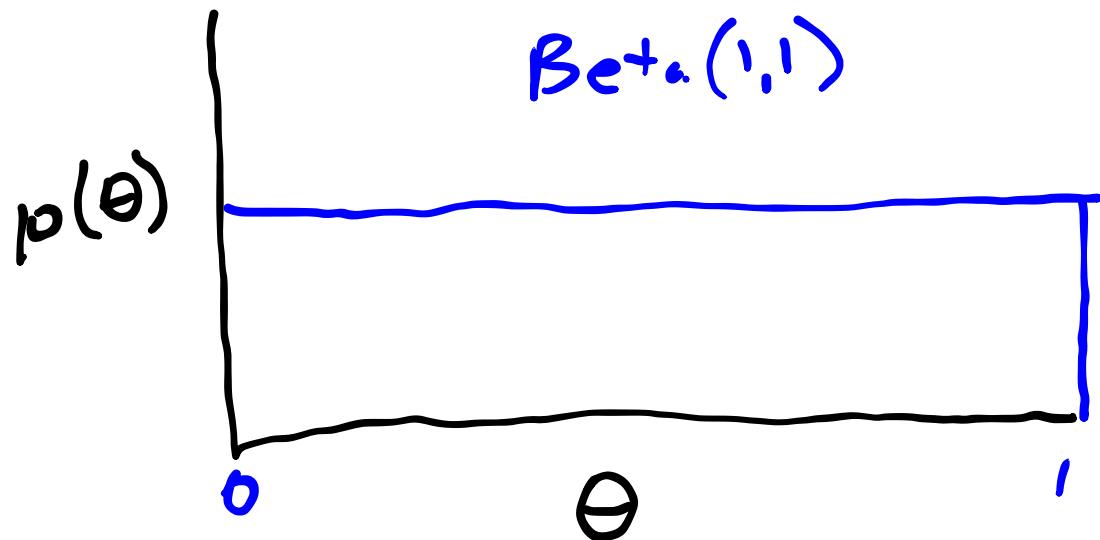
Given a $\text{Beta}(1, 1)$ prior distribution

The posterior distribution of θ is
 $\text{Beta}(w + 1, l + 1)$

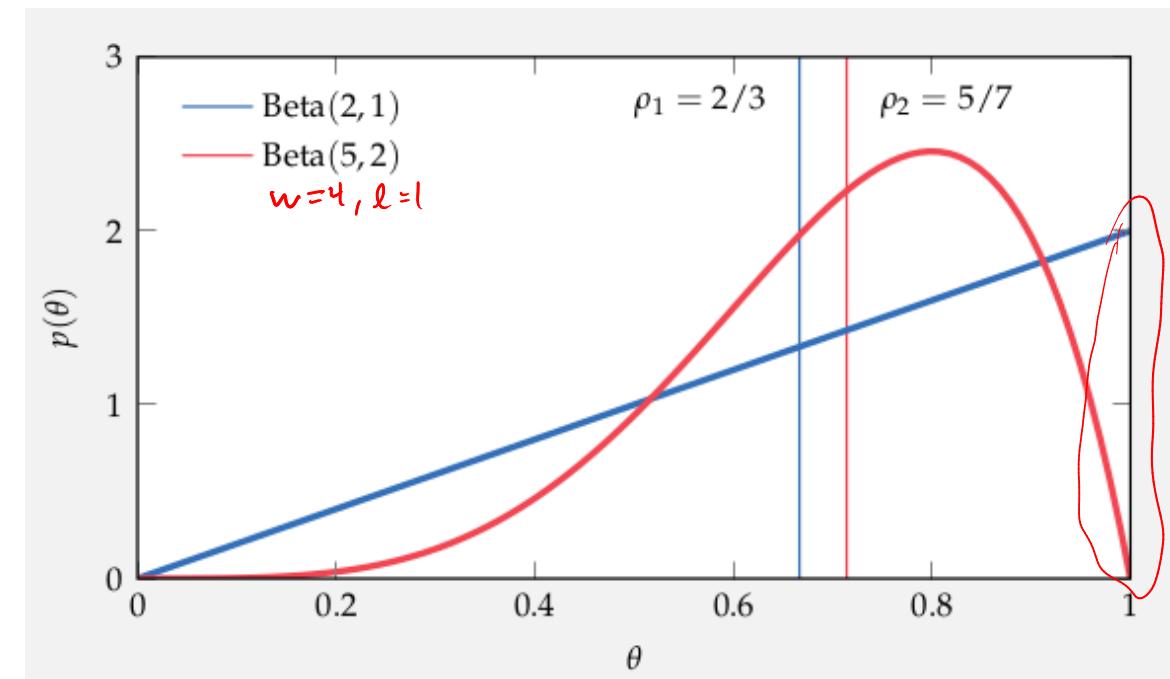


Bayesian Estimation

Given a $\text{Beta}(1, 1)$ prior distribution



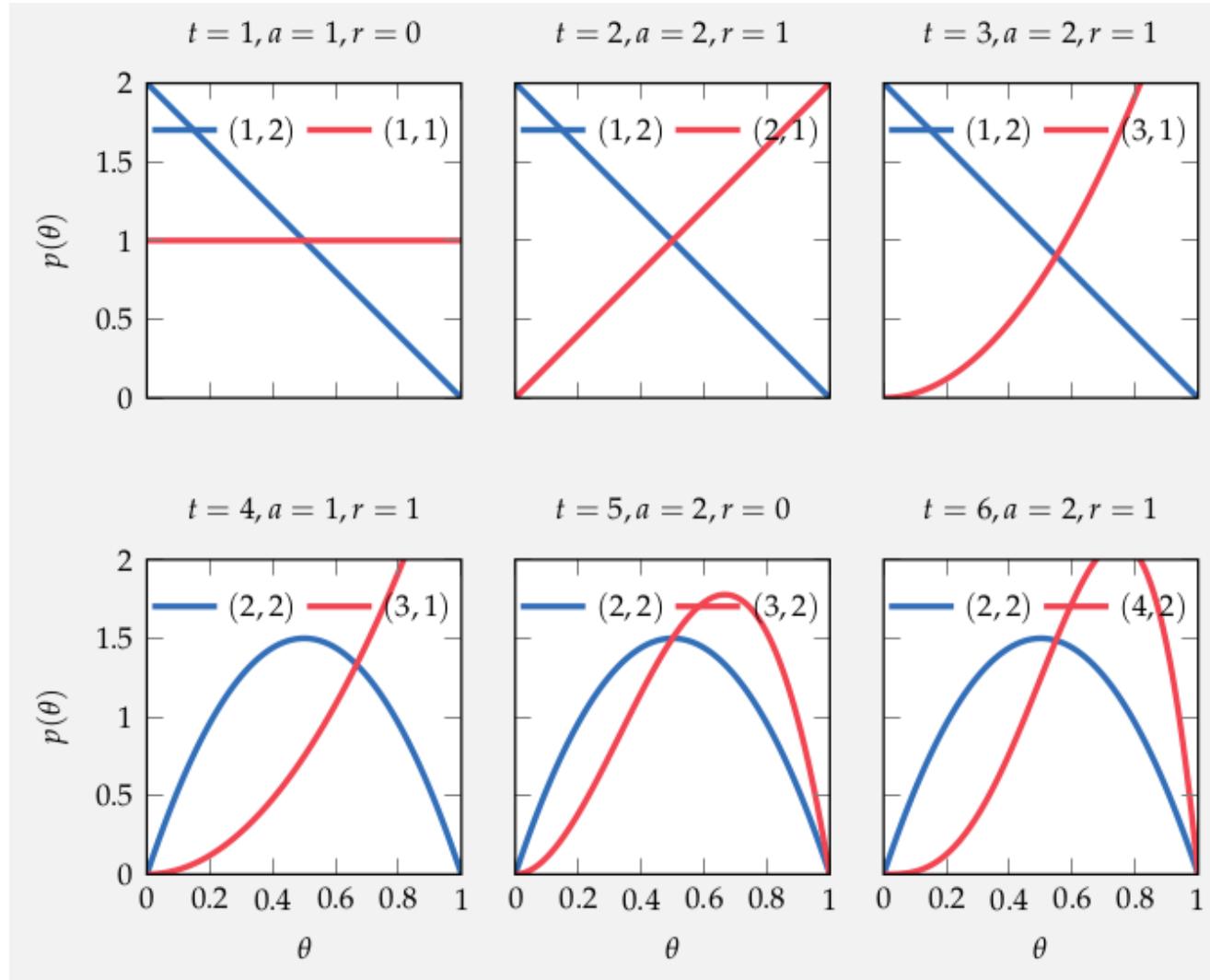
The posterior distribution of θ is
 $\text{Beta}(w + 1, l + 1)$



Bayesian Estimation

arm 1

arm 2



Bayesian Bandit Algorithms

Bayesian Bandit Algorithms

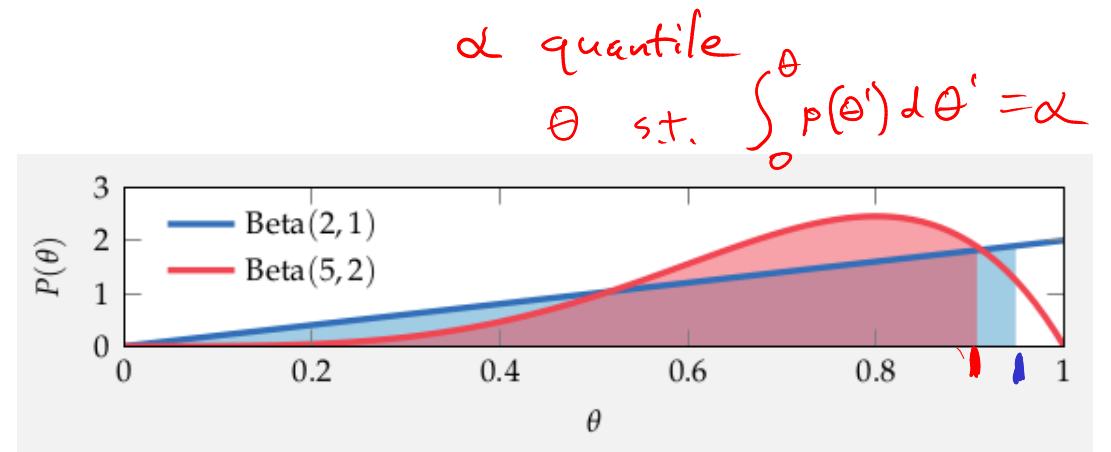
- Quantile Selection

Choose a for which the α quantile of
 $b(\theta)$ is highest

Bayesian Bandit Algorithms

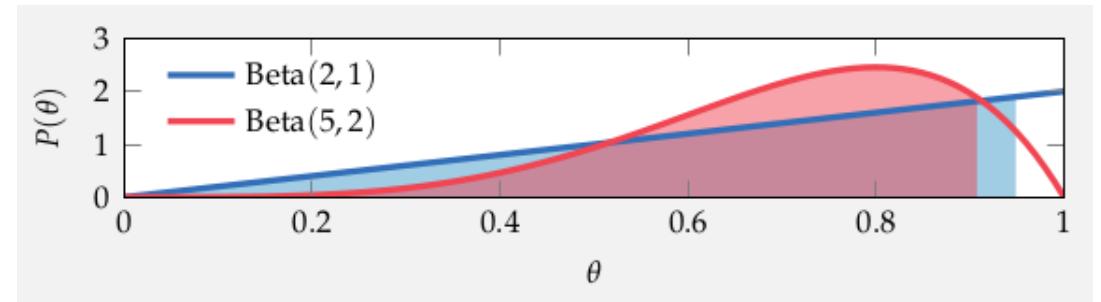
- Quantile Selection

Choose a for which the α quantile of $b(\theta)$ is highest

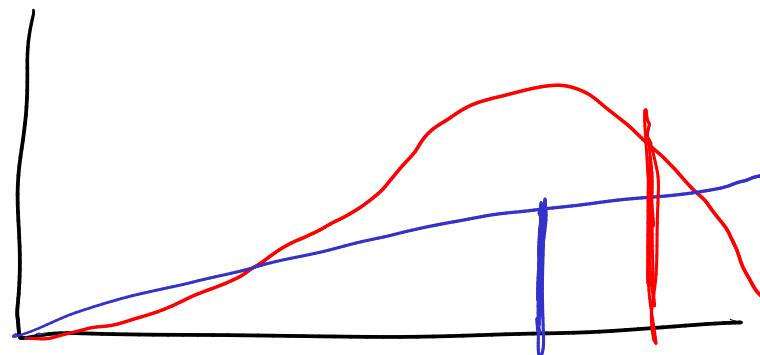


Bayesian Bandit Algorithms

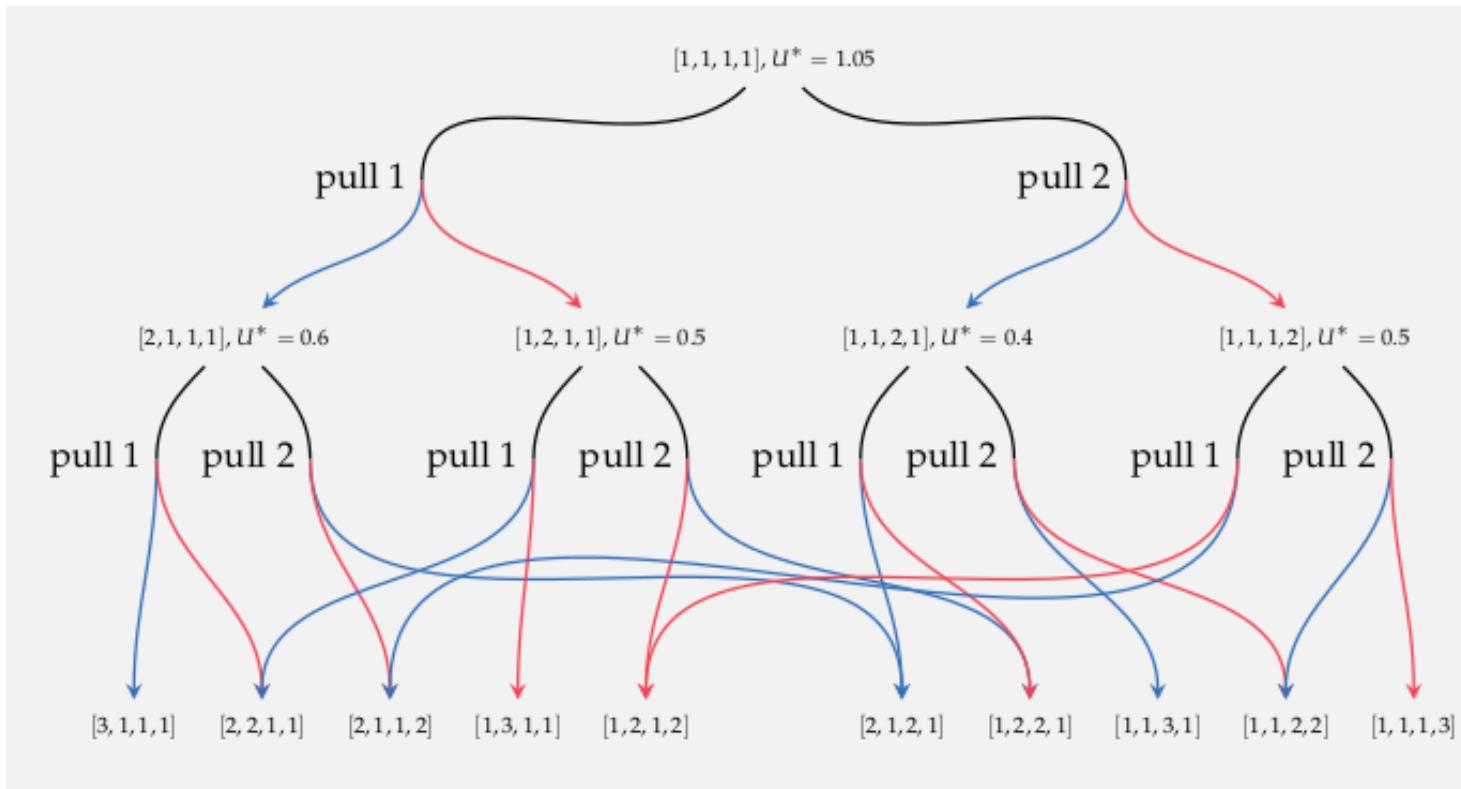
- Quantile Selection
Choose a for which the α quantile of $b(\theta)$ is highest



- Thompson Sampling
Sample $\hat{\theta}$
Choose $\operatorname{argmax}_a \hat{\theta}_a$



Optimal Algorithm - Dynamic Programming



Easier to Implement
Faster

- ↑ - greedy
- ϵ -greedy
- explore-commit
- softmax

$$\alpha e^{\lambda p_a}$$

Review

Optimal in Limit

$$\begin{aligned}\epsilon &\rightarrow 0 \\ k &\rightarrow \infty \\ \lambda &\rightarrow \infty\end{aligned}$$

= UCB

$$p_a + C \sqrt{\frac{\log N}{N(a)}}$$



- Quantile Selection



- Thompson Sampling



- Dynamic Programming



Bayesian

Less Regret

for a Bernoulli Bandit
Regret $\equiv \Theta^* N - \sum_{t=1}^N r_t$
 $O(N)$

$O(N)$

$O(N)$

$O(N)$

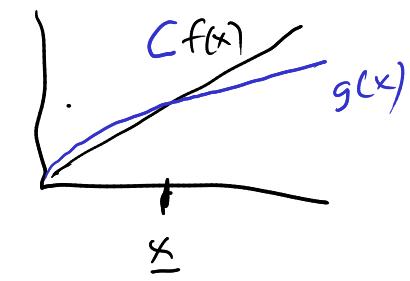
$O(\log(N))$

$g(x) = O(f(x))$

$\exists C$ s.t.

$g(x) \leq C f(x)$

$\forall x > x_*$



Guiding Questions

- What are the best ways to trade off Exploration and Exploitation