

( $s, a, r, \gamma$ )  
??

# Last Time

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

# Last Time

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

# Last Time

# Last Time

First RL Algorithm:

# Last Time

First RL Algorithm:

Tabular Maximum Likelihood Model-Based Reinforcement Learning

# Last Time

First RL Algorithm:

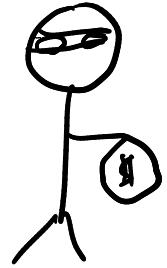
Tabular Maximum Likelihood Model-Based Reinforcement Learning

loop  
choose action  $a$    $\epsilon/\bar{\epsilon}$   
gain experience  
estimate  $T, R$   
solve MDP with  $T, R$

# Guiding Questions

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

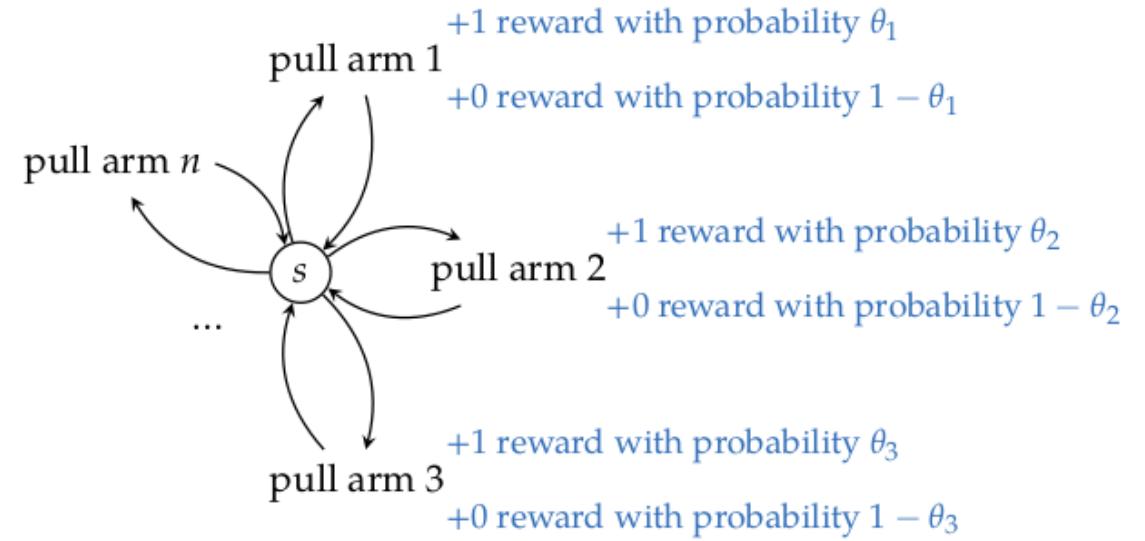
# Bandits



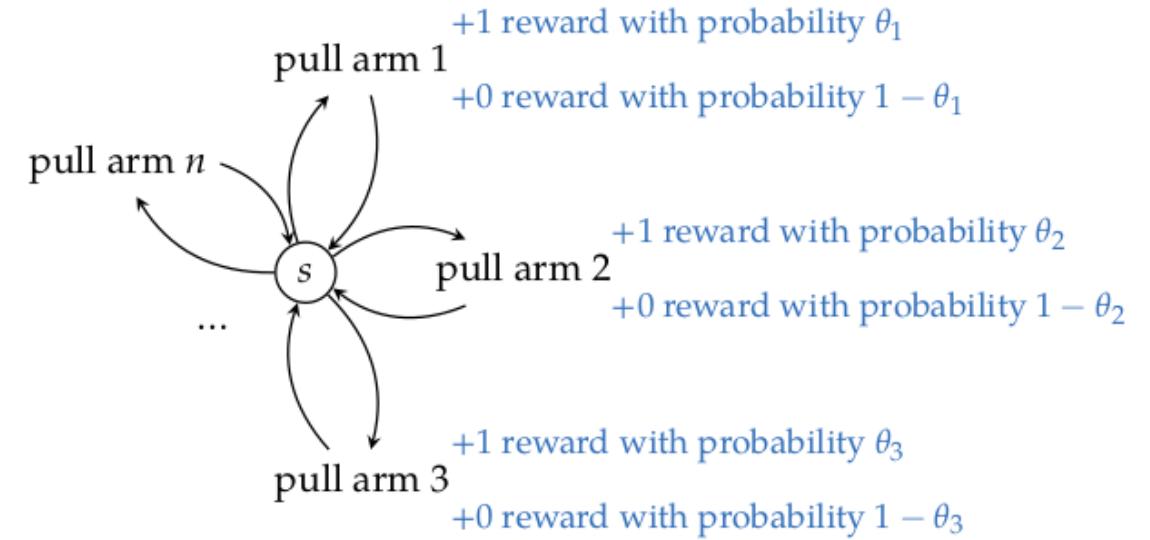
# Bandits



# Bandits

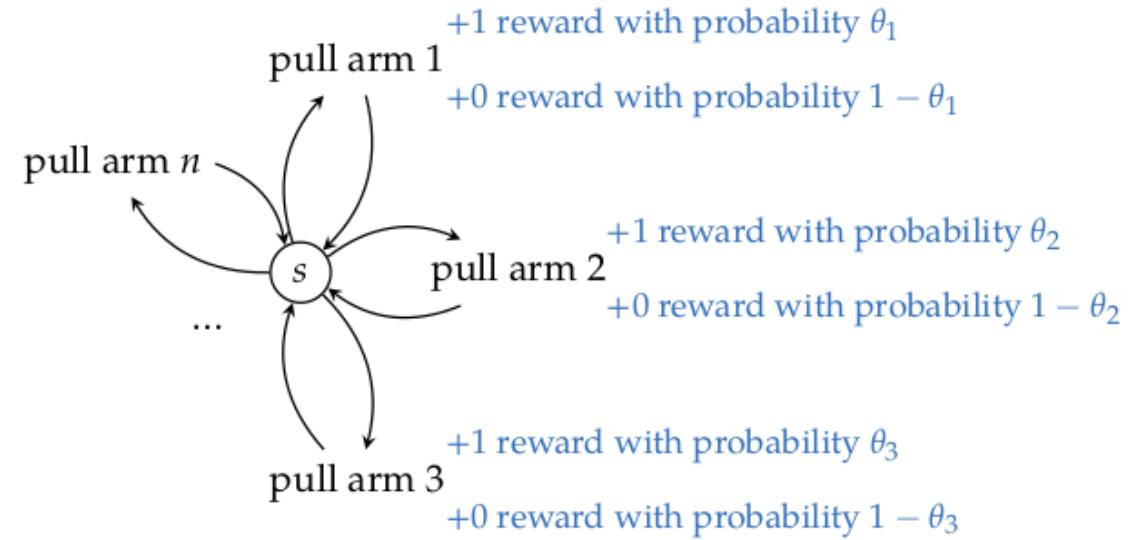


# Bandits



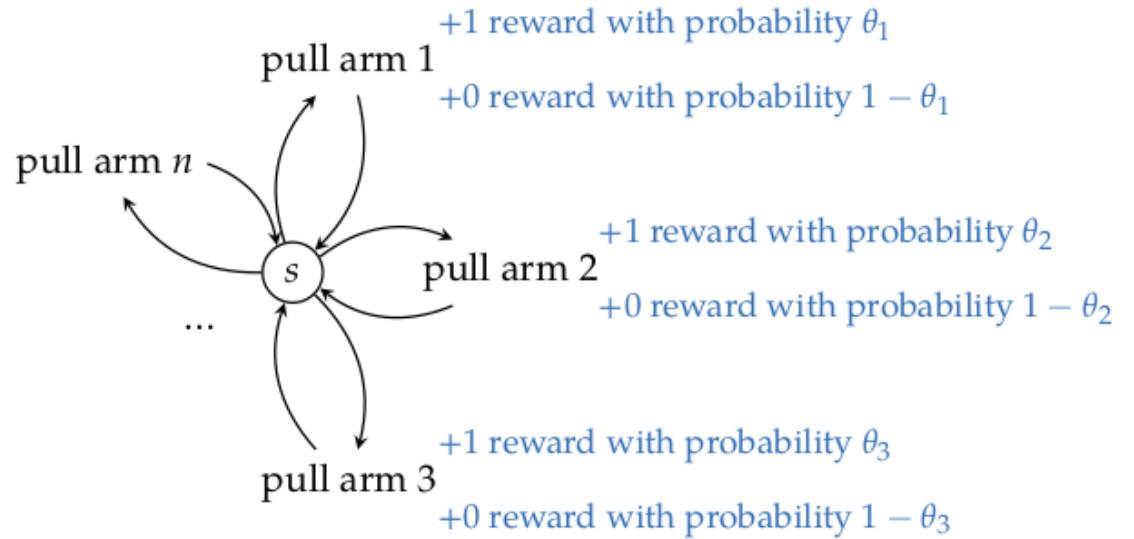
- Bernoulli Bandit with parameters  $\theta$

# Bandits



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

# Bandits



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

“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

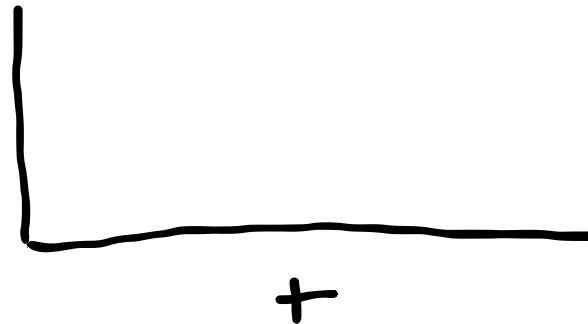
- Explore then Commit

Choose  $a$  randomly for  $k$  steps

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

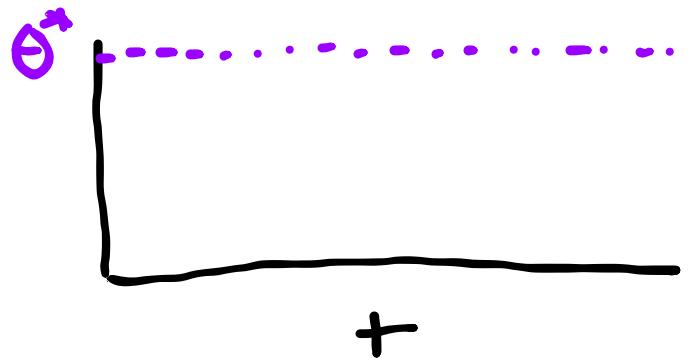
# Undirected Strategies

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



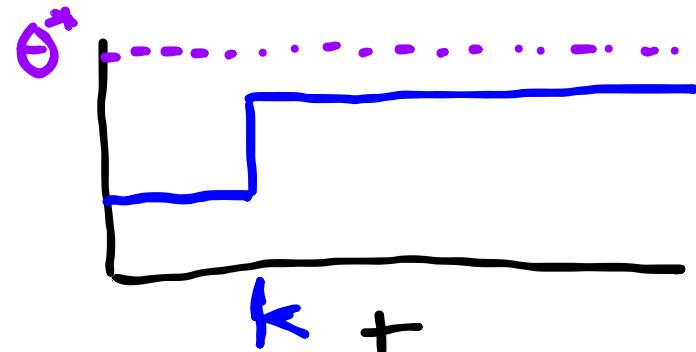
# Undirected Strategies

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



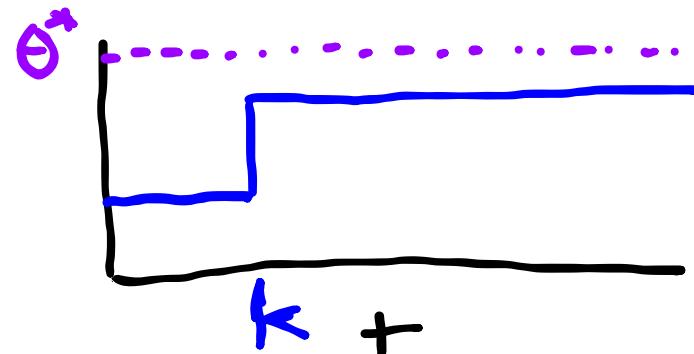
# Undirected Strategies

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



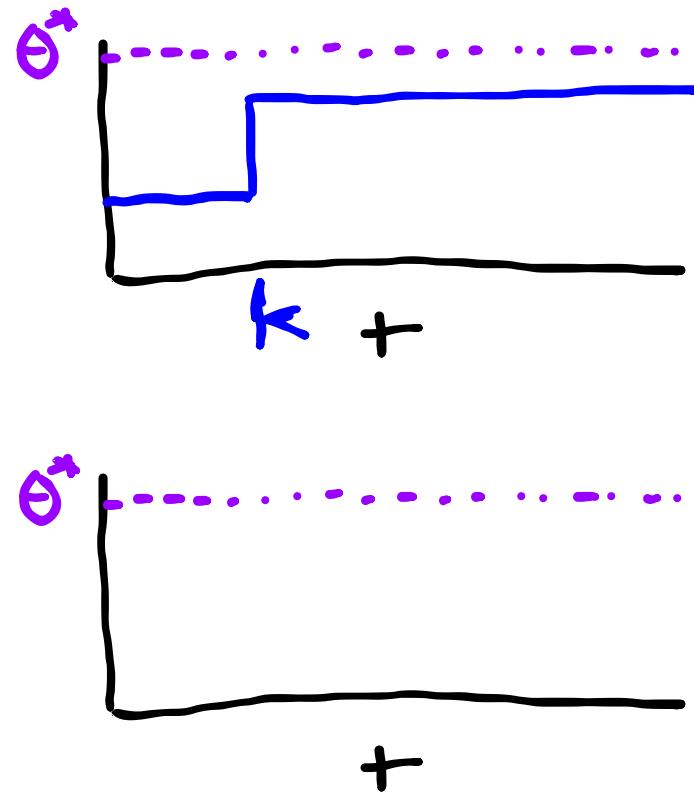
# Undirected Strategies

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



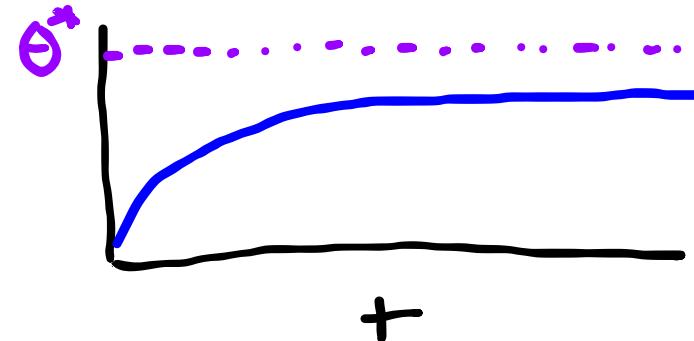
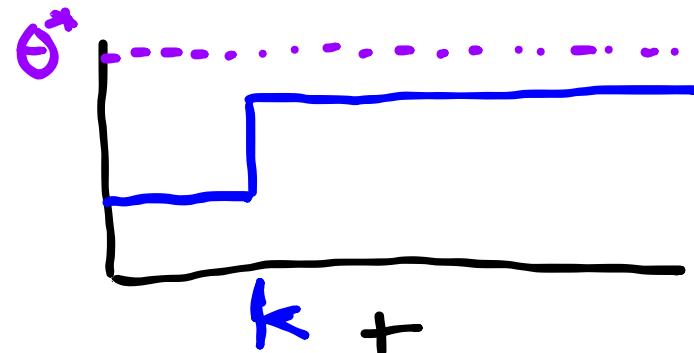
# Undirected Strategies

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



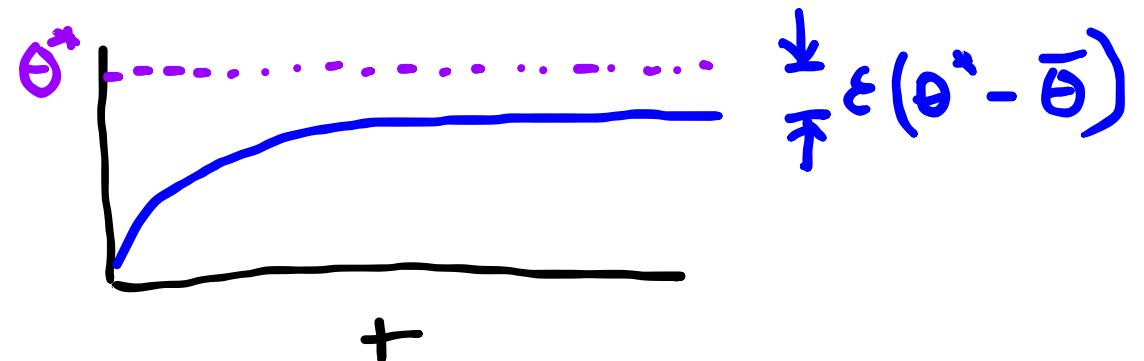
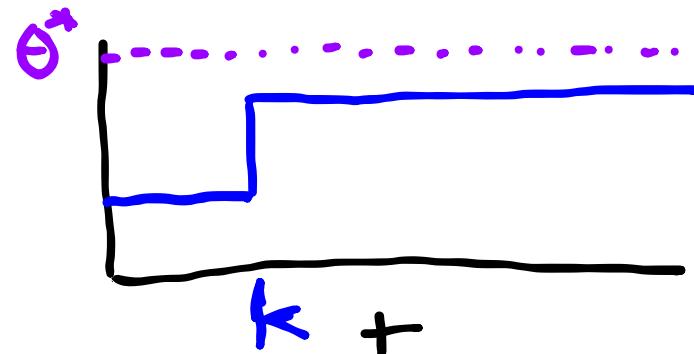
# Undirected Strategies

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



# Undirected Strategies

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



# Directed Strategies

(Non-Bayesian)



# Directed Strategies

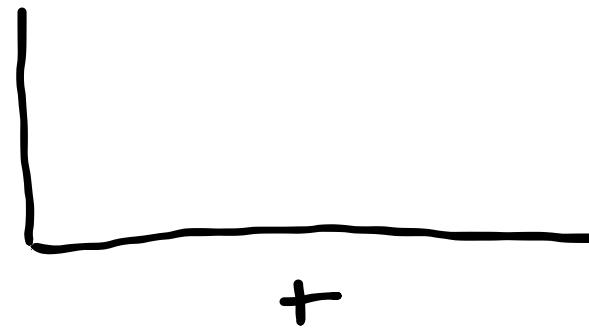
- Softmax

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



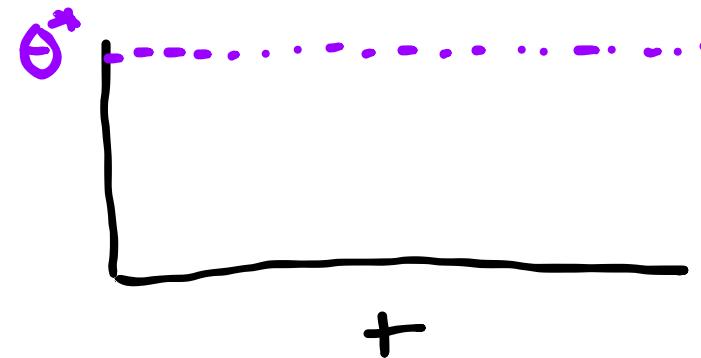
# Directed Strategies

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



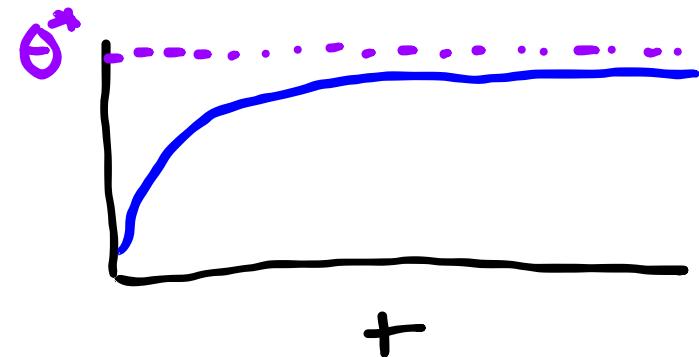
# Directed Strategies

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



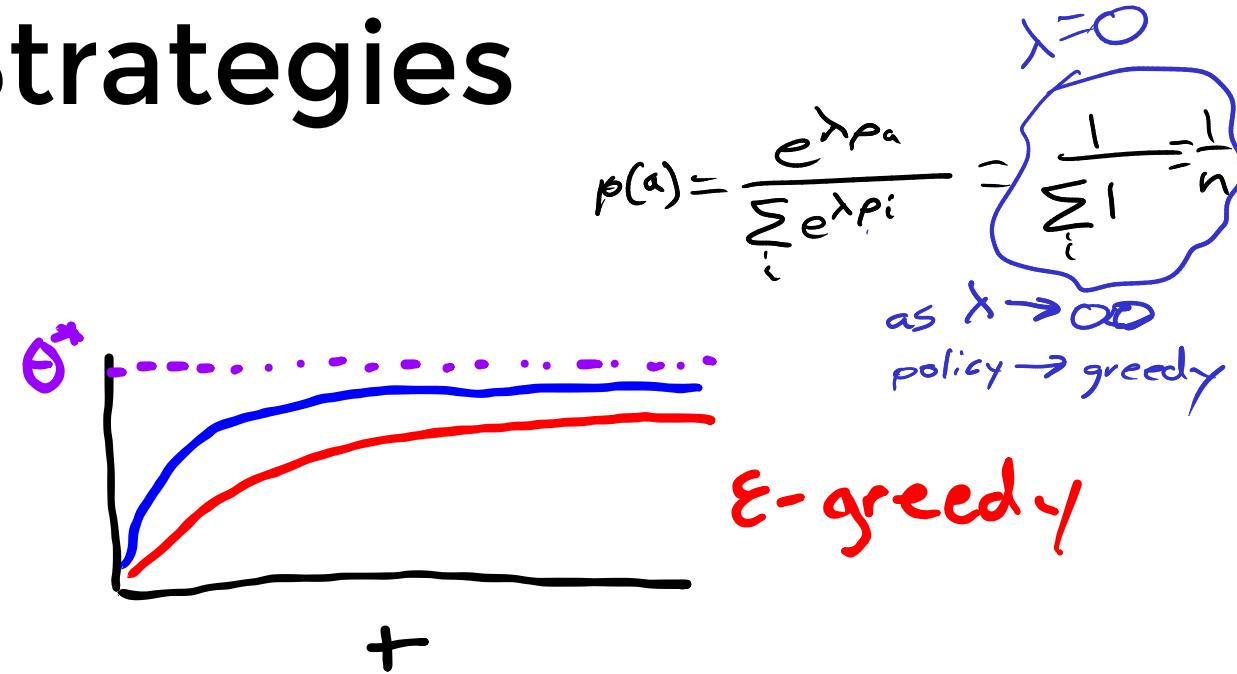
# Directed Strategies

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



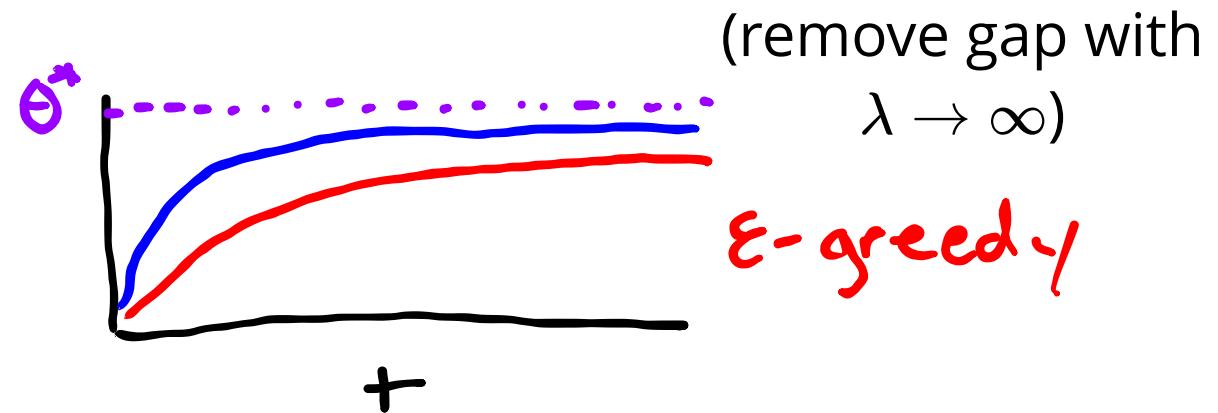
# Directed Strategies

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



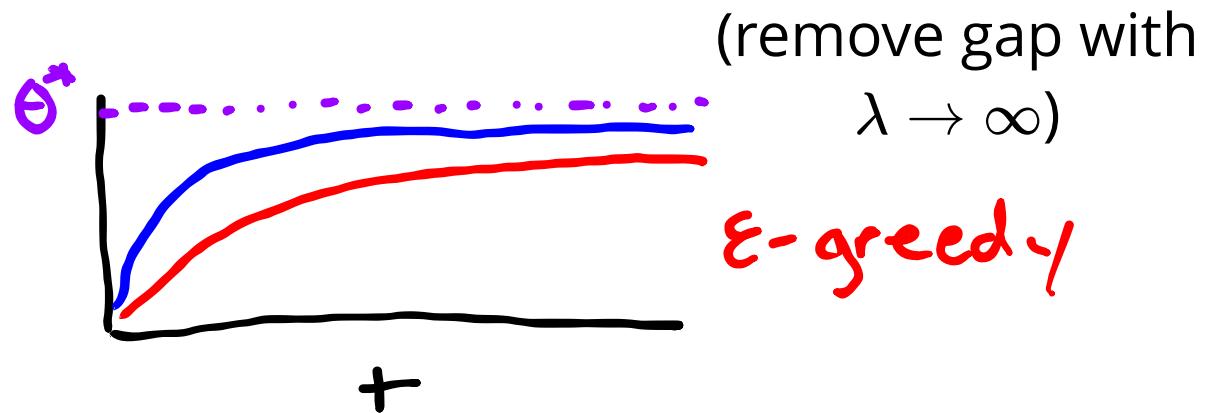
# Directed Strategies

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



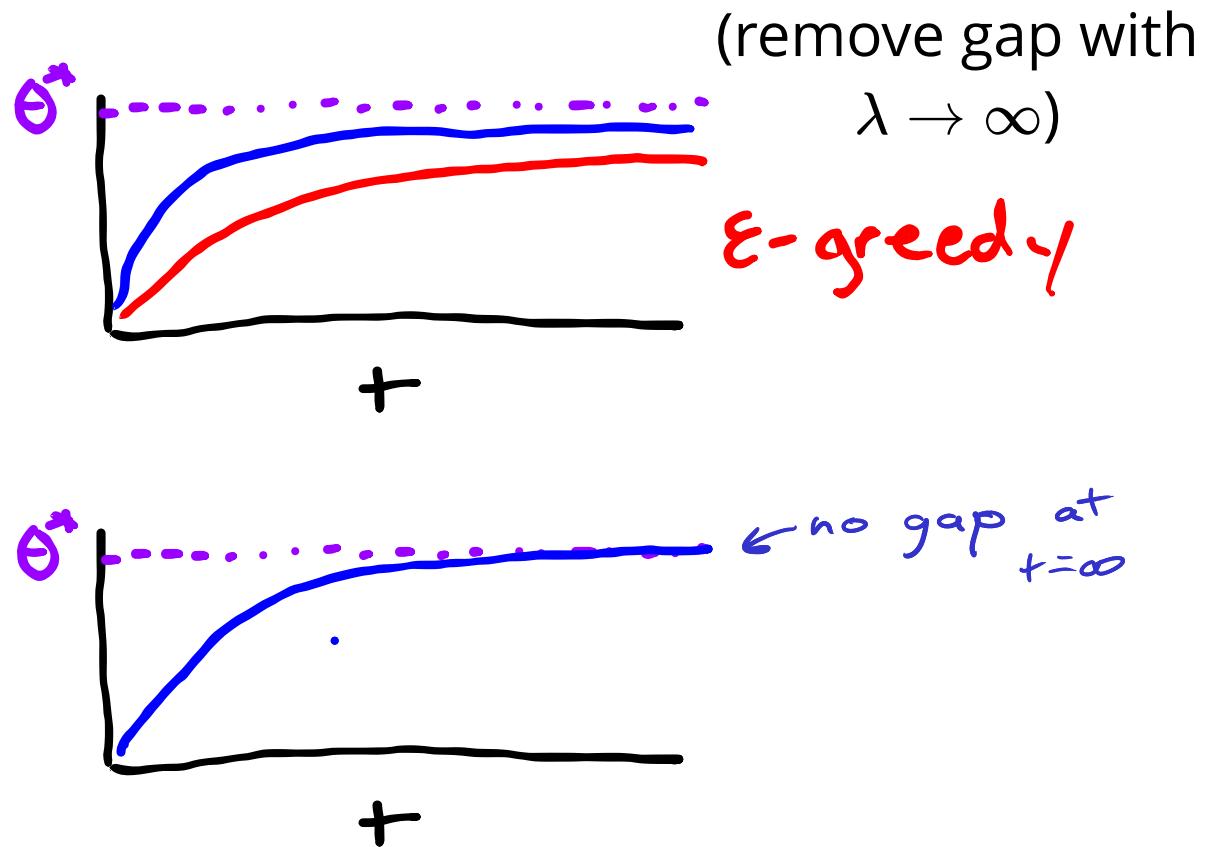
# Directed Strategies

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



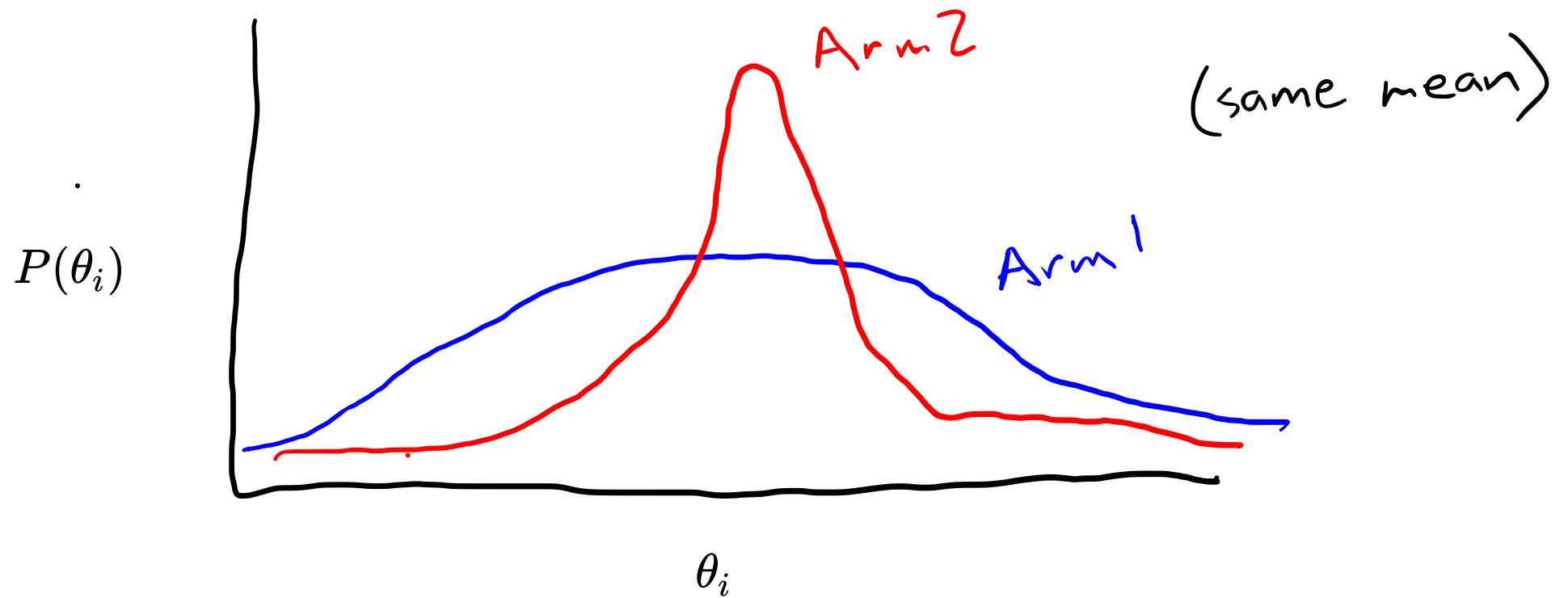
# Directed Strategies

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



# Break

Discuss with your neighbor: Suppose you have the following *belief* about the parameters  $\theta$ . Which arm should you choose to pull next?



# Bayesian Estimation

# Bayesian Estimation

Bernoulli Distribution

# Bayesian Estimation

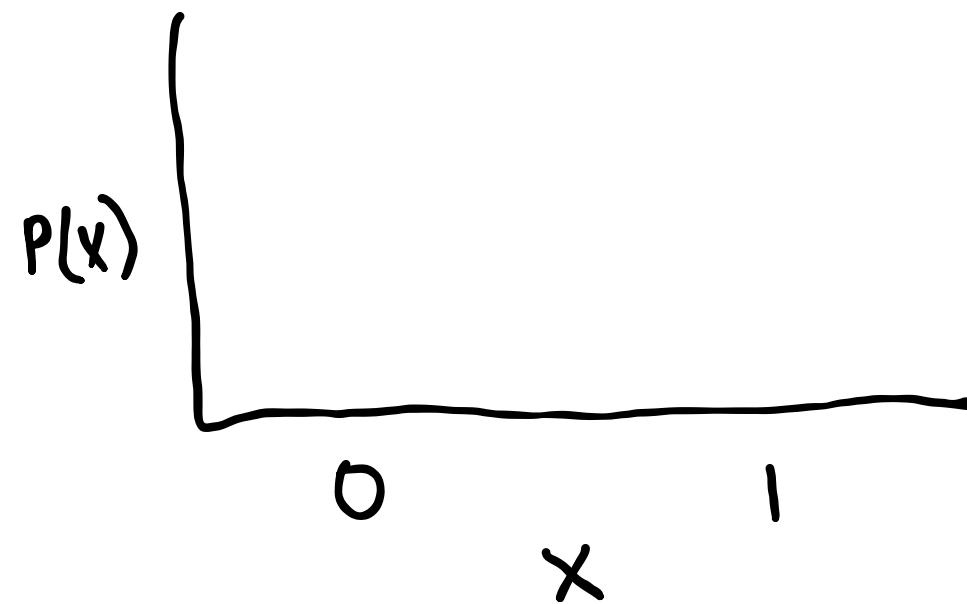
Bernoulli Distribution

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

# Bayesian Estimation

Bernoulli Distribution

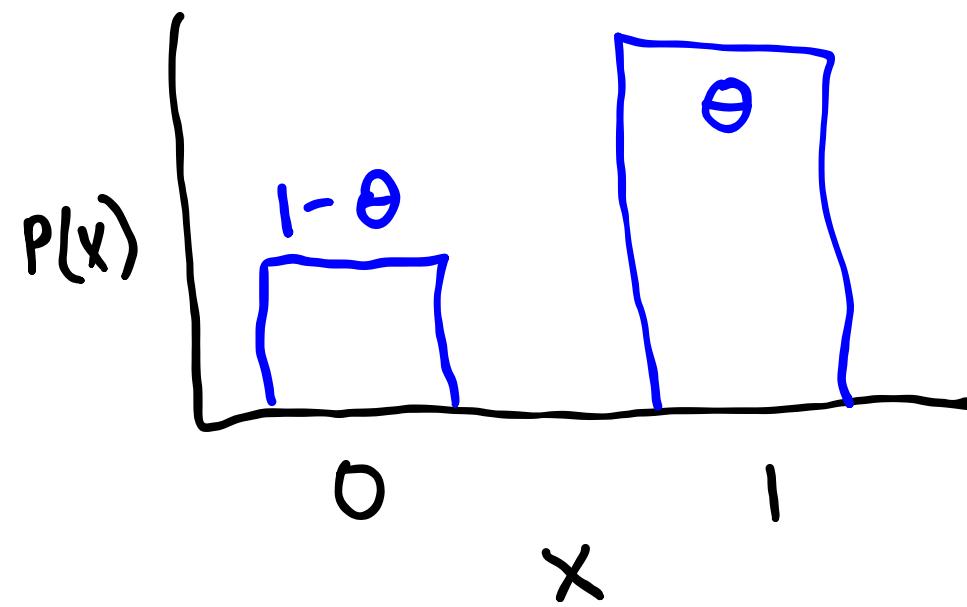
$\text{Bernoulli}(\theta)$



# Bayesian Estimation

Bernoulli Distribution

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

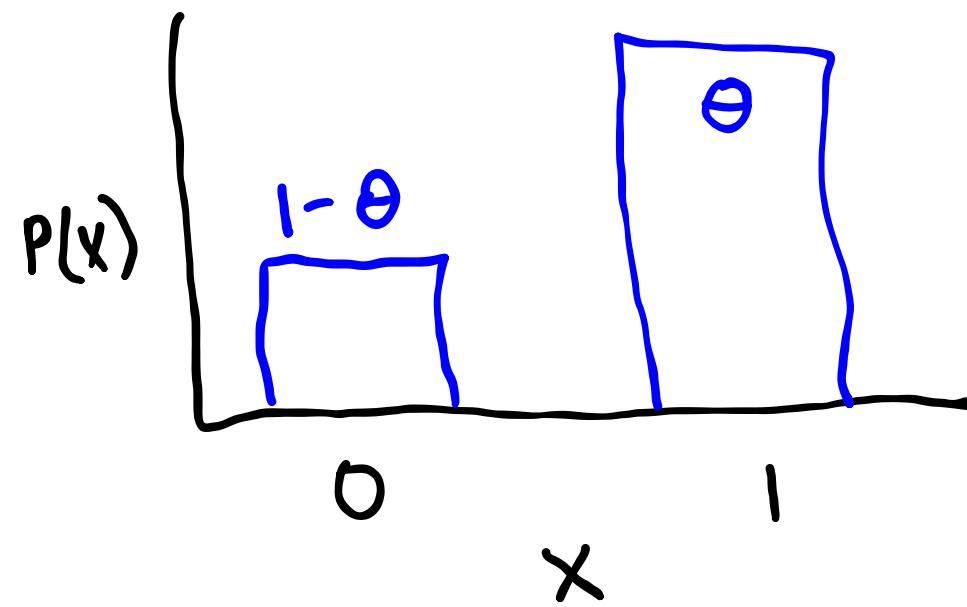


# Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

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

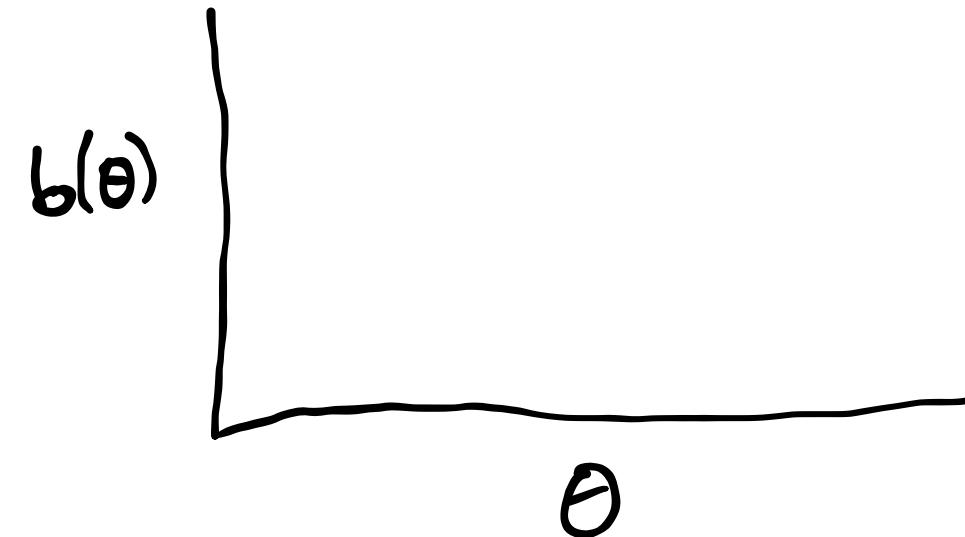
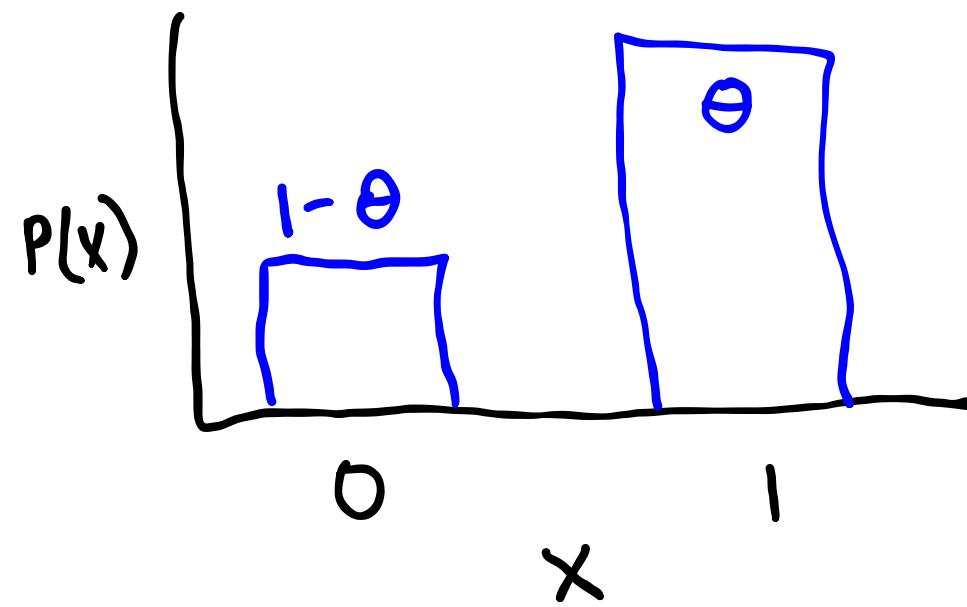


# Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

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

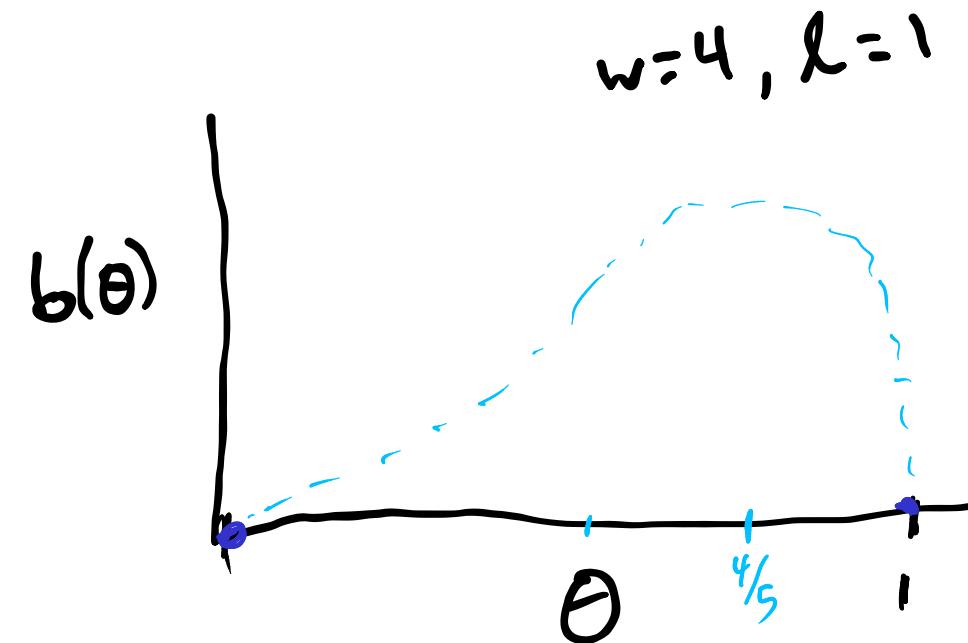
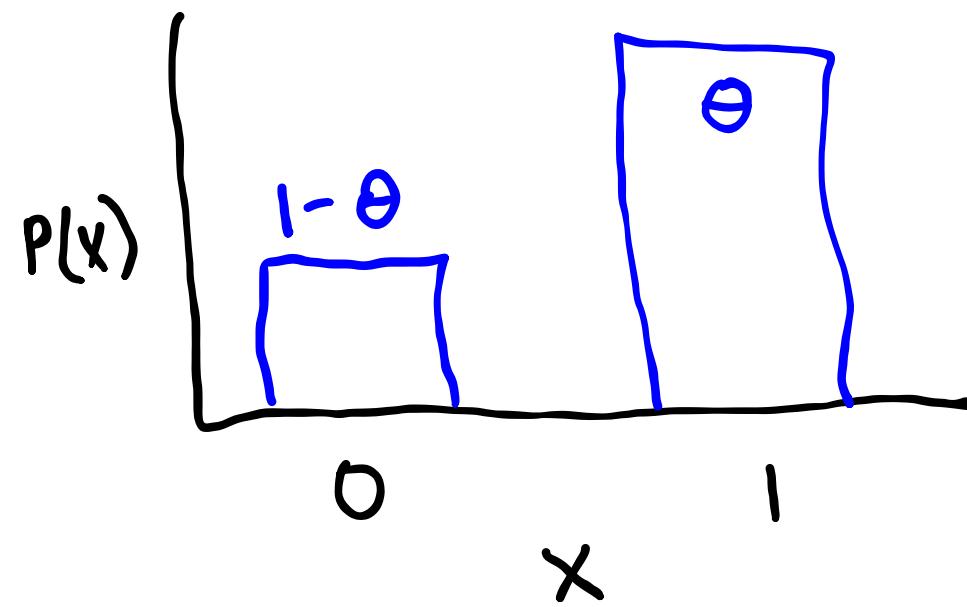


# Bayesian Estimation

Bernoulli Distribution

$\text{Bernoulli}(\theta)$

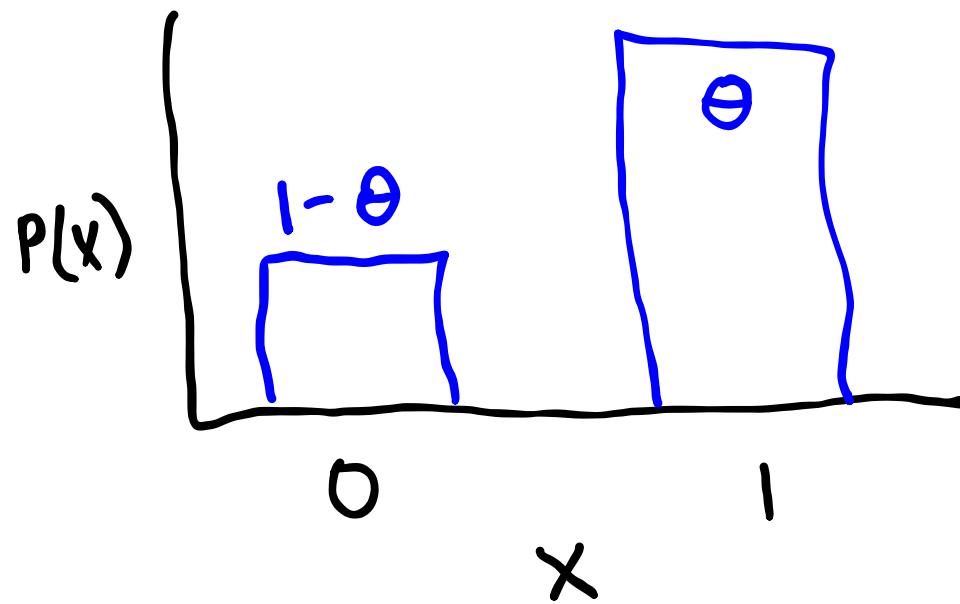
Discussion: Given that I have received  $w$  wins and  $l$  losses, what should my belief (probability distribution) about  $\theta$  look like?



# Bayesian Estimation

Bernoulli Distribution

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



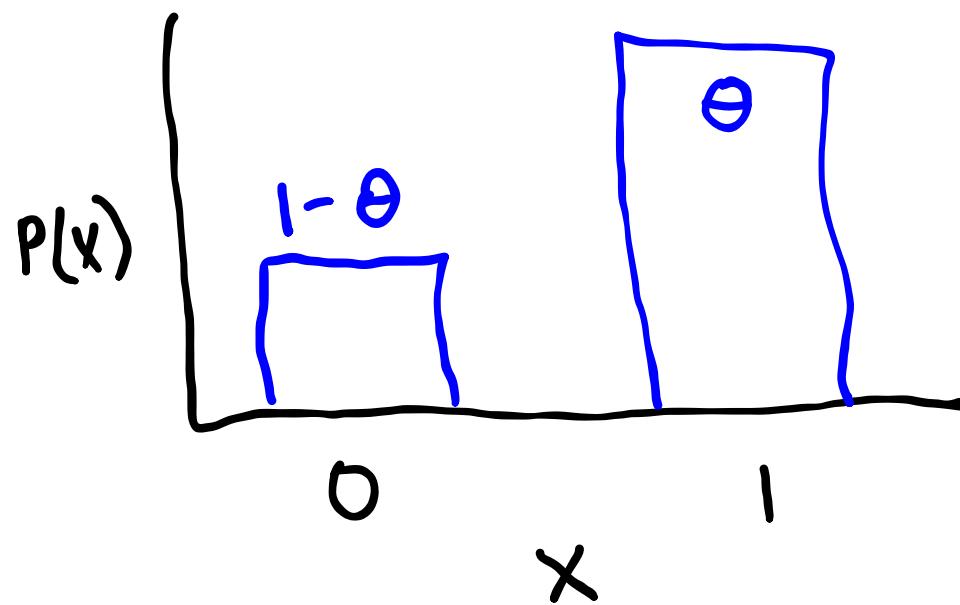
# Bayesian Estimation

Bernoulli Distribution

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

Beta Distribution

(distribution over Bernoulli distributions)



# Bayesian Estimation

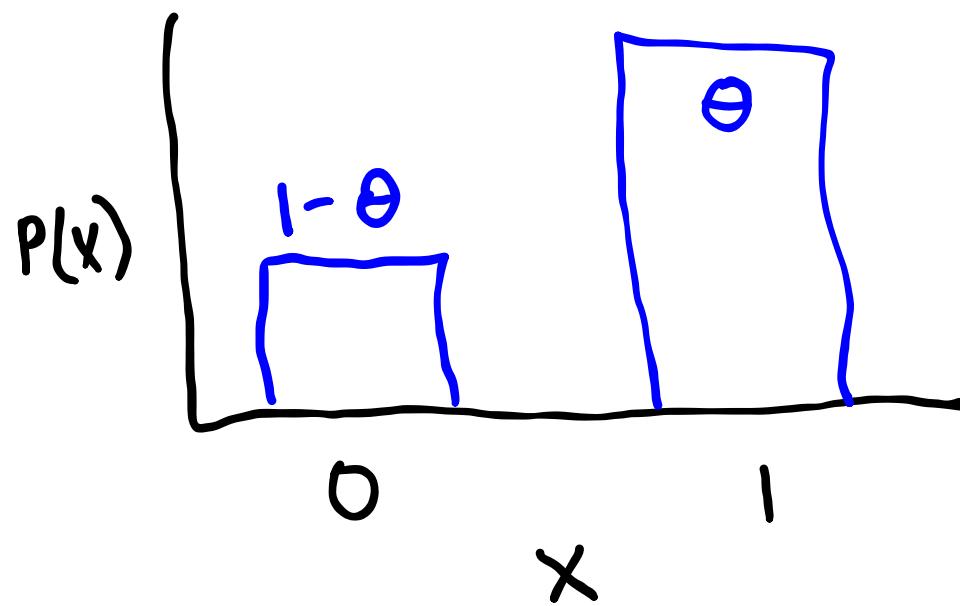
Bernoulli Distribution

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

Beta Distribution

(distribution over Bernoulli distributions)

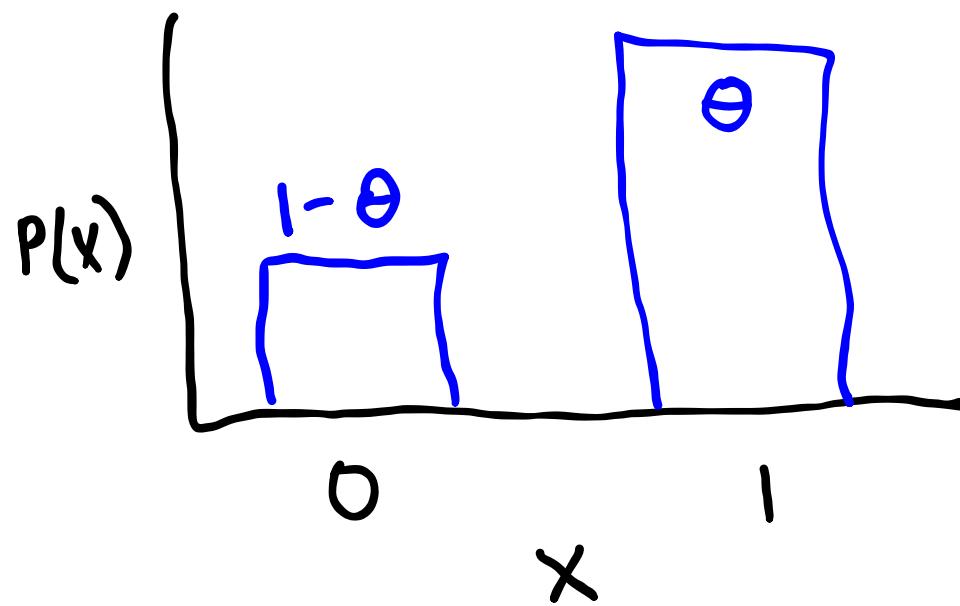
$$\text{Beta}(\alpha, \beta)$$



# Bayesian Estimation

Bernoulli Distribution

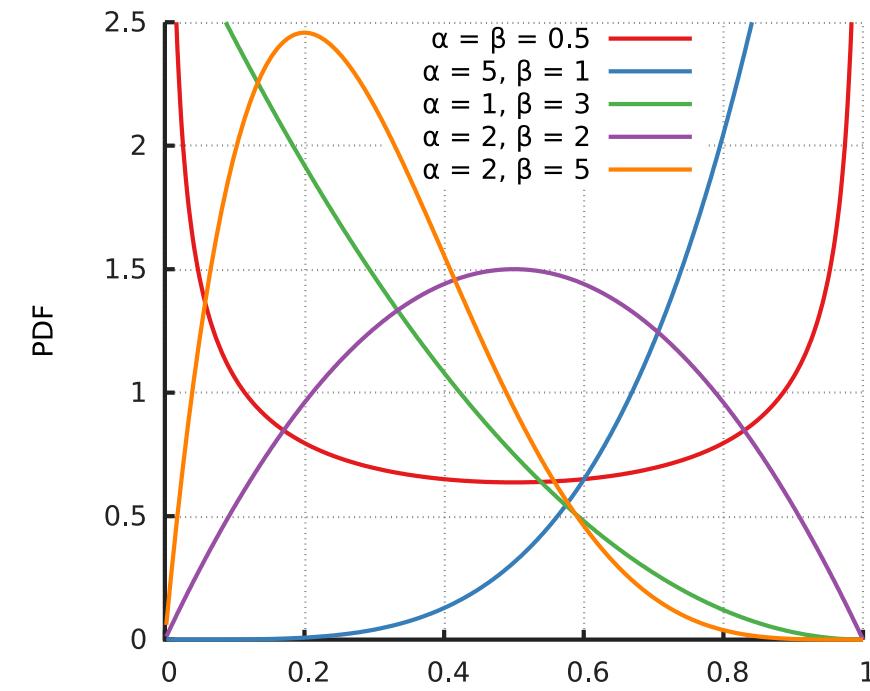
$\text{Bernoulli}(\theta)$



Beta Distribution

(distribution over Bernoulli distributions)

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



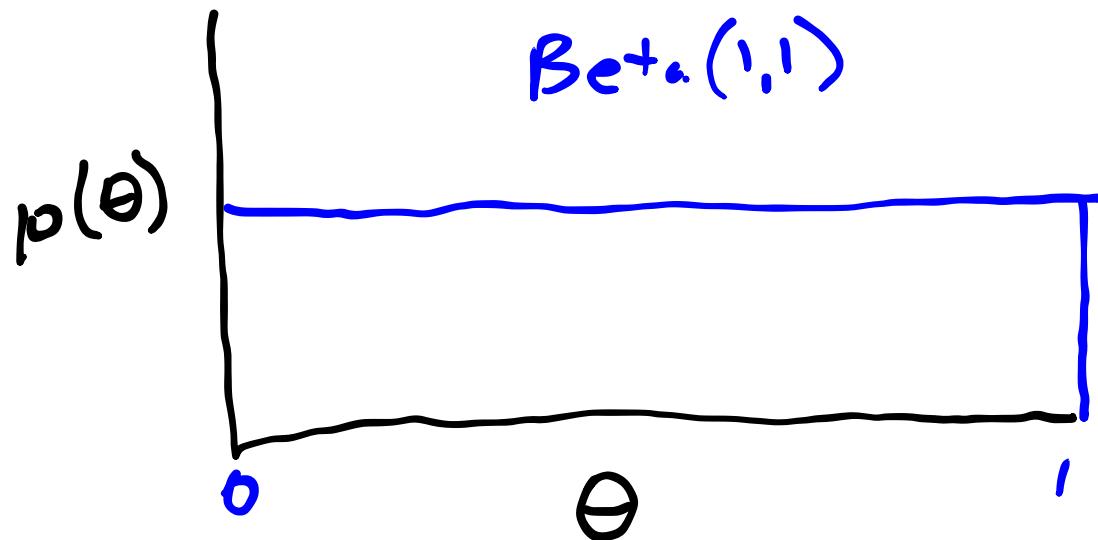
# Bayesian Estimation

# Bayesian Estimation

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

# Bayesian Estimation

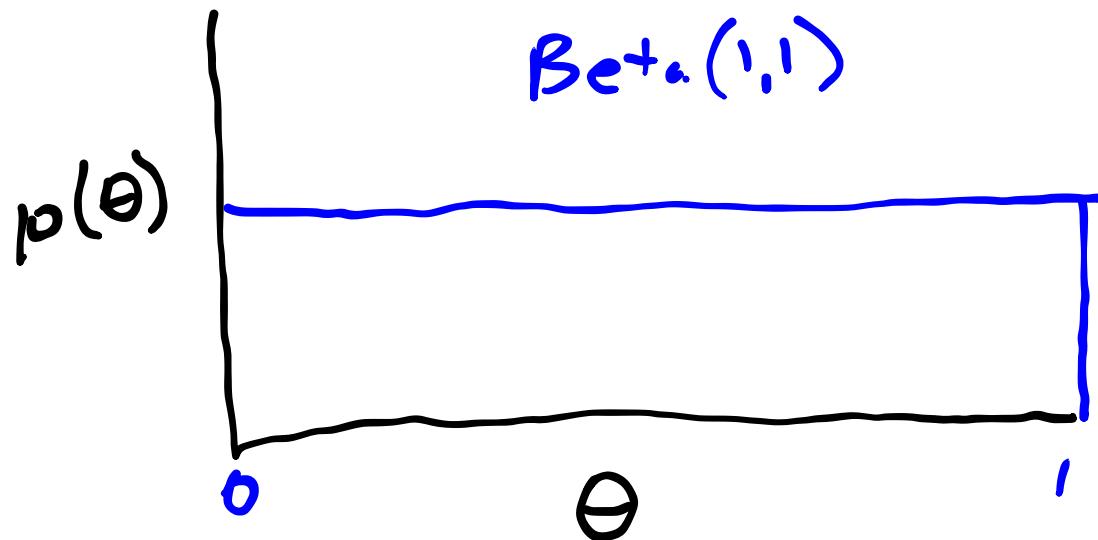
Given a  $\text{Beta}(1, 1)$  prior distribution



# Bayesian Estimation

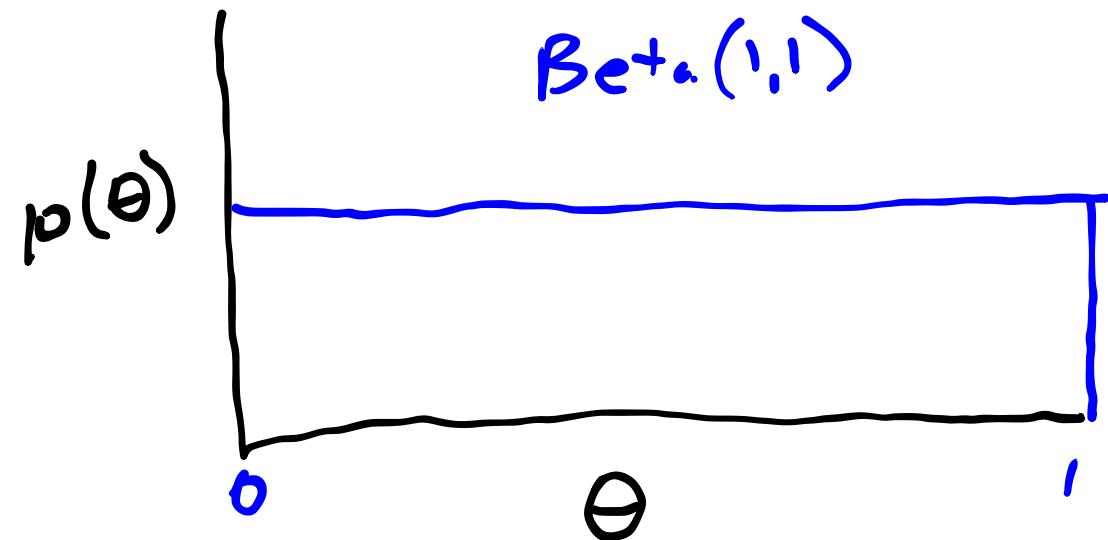
Given a  $\text{Beta}(1, 1)$  prior distribution

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

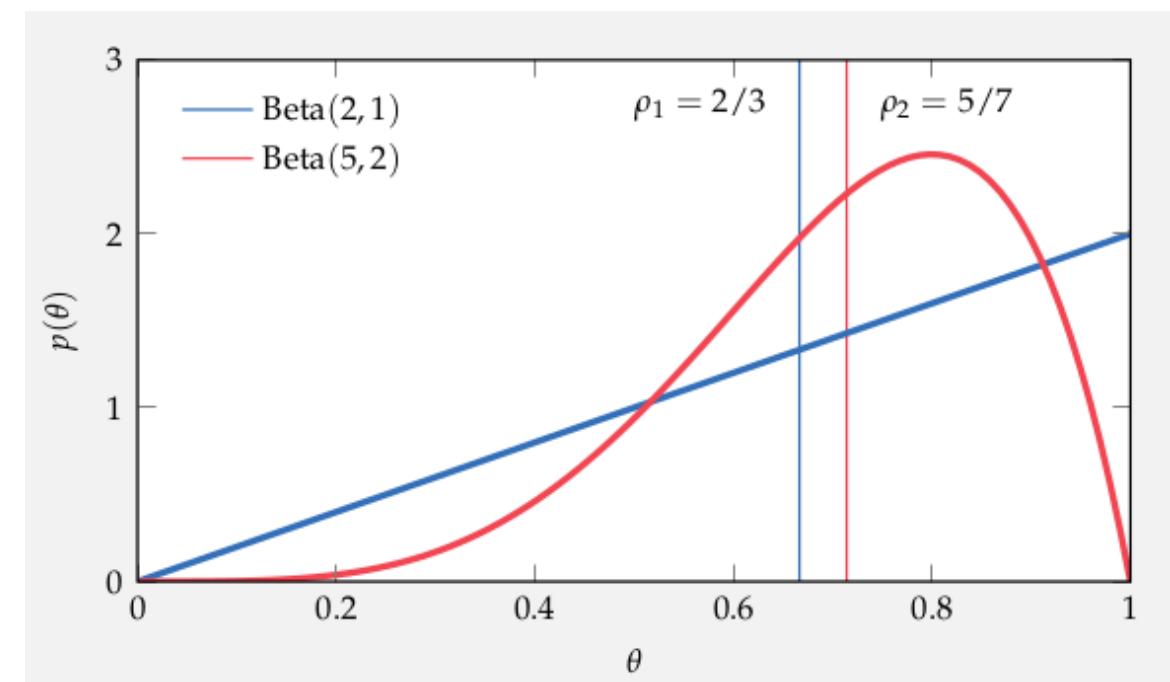


# Bayesian Estimation

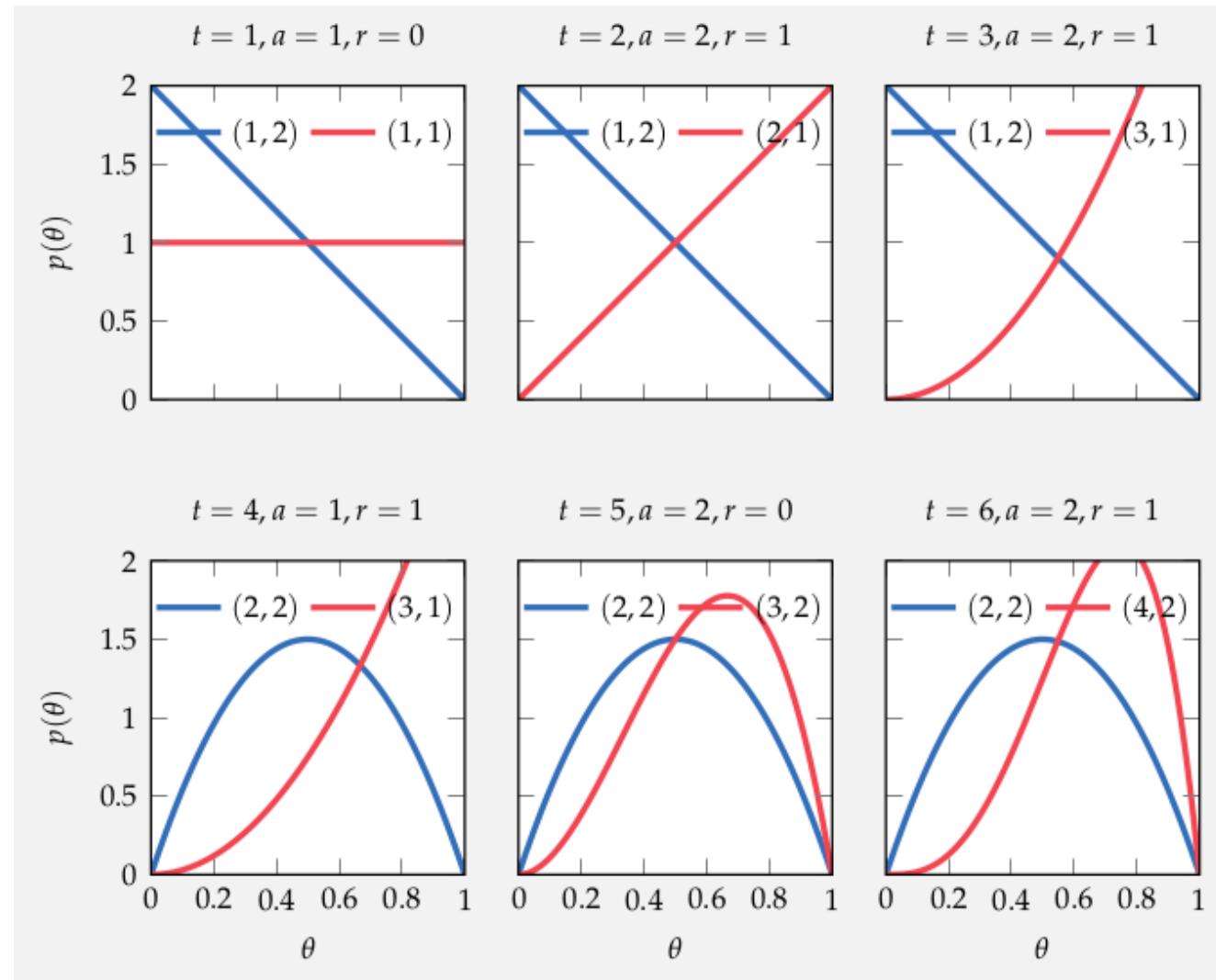
Given a  $\text{Beta}(1, 1)$  prior distribution



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



# Bayesian Estimation



$t$  = time

$a$  = arm pulled

$r$  = reward

# Bayesian Bandit Algorithms

# Bayesian Bandit Algorithms

- Quantile Selection

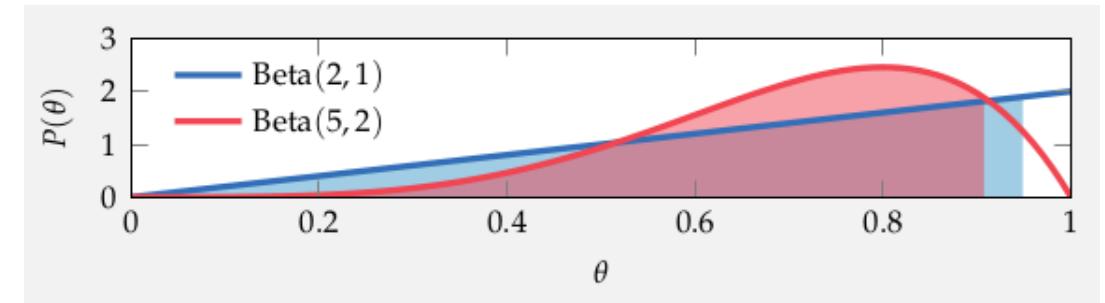
Choose  $a$  for which the  $\alpha$  quantile of  
 $b(\theta)$  is highest

# Bayesian Bandit Algorithms

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

$\alpha = 0.9$

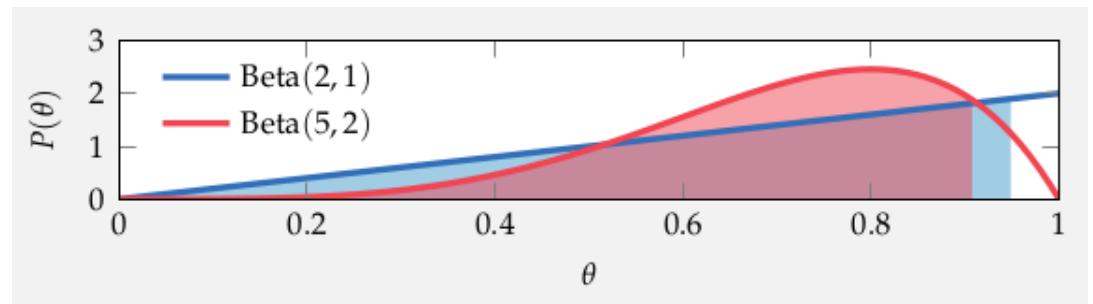
higher  $\alpha$   
more optimistic



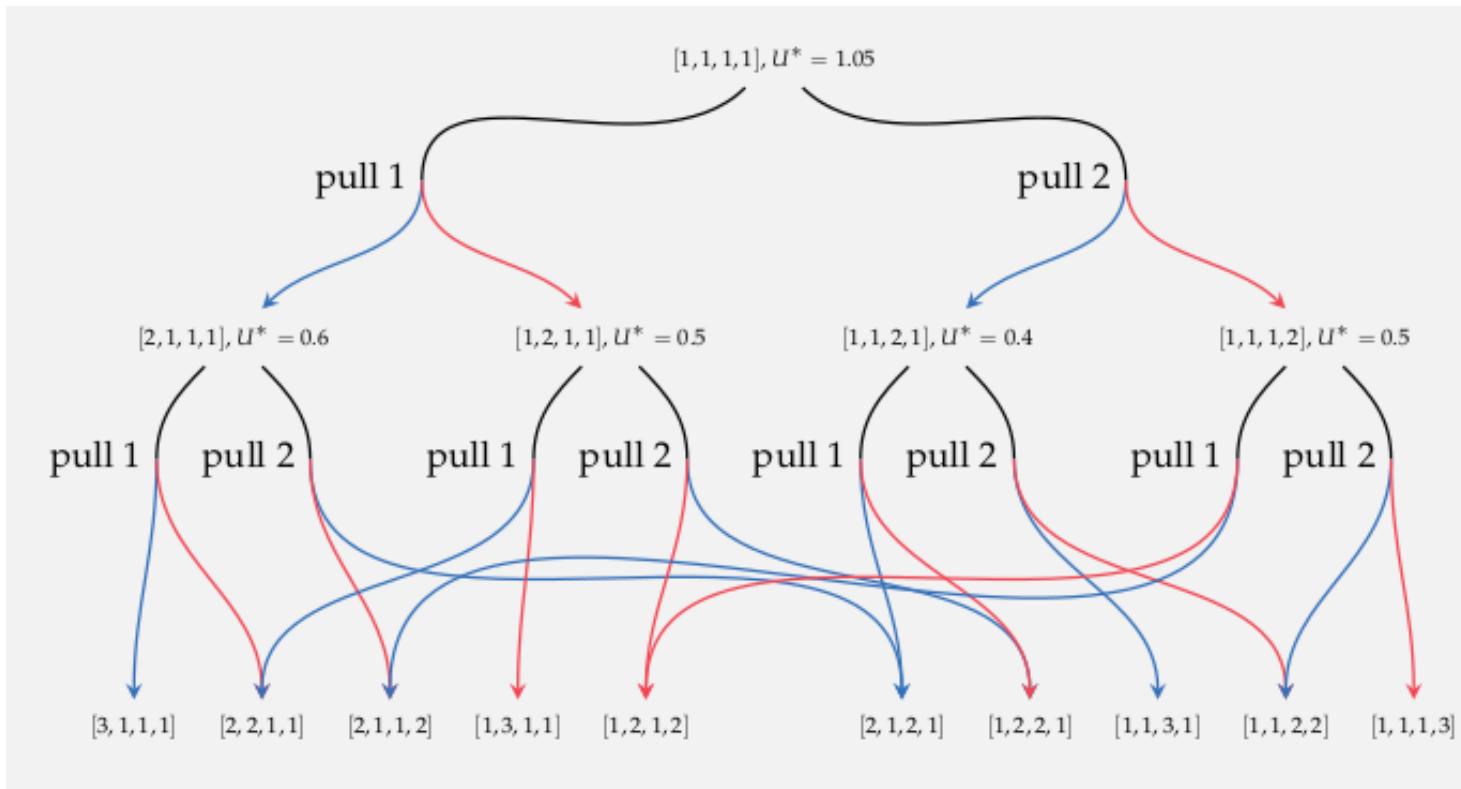
# Bayesian Bandit Algorithms

- Quantile Selection  
Choose  $a$  for which the  $\alpha$  quantile of  $b(\theta)$  is highest
- Thompson Sampling  
Sample  $\hat{\theta}$   
Choose  $\operatorname{argmax}_a \hat{\theta}_a$

$$\alpha = 0.9$$



# Optimal Algorithm - Dynamic Programming



Easier  
to Implement,  
Faster



# Review

Optimal in limit

for a Bernoulli  
Bandit

$$\text{Regret} \equiv \theta^* N - \sum_{t=1}^N r_t$$

- Undirected
- greedy
  - $\epsilon$  greedy
  - explore-commit
  - softmax

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

No

$$\epsilon \rightarrow 0$$

$$k \rightarrow \infty$$

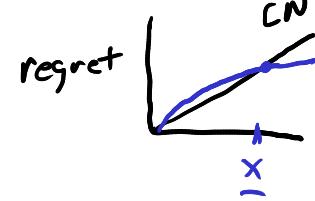
$$\lambda \rightarrow \infty$$

$$O(N)$$

$$O(N)$$

$$O(N)$$

$$O(N)$$



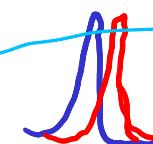
- UCB

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



$$O(\log(N))$$

- Bayesian
- Quantile Selection
  - Thompson Sampling
  - Dynamic Programming



"



"



$$g(x) = O(f(x))$$
$$\exists C \in \mathbb{R}^+$$
$$g(x) \leq C f(x)$$
$$x > 7$$

Less Regret

# Guiding Questions

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

UCB