

20. Topics in Online Learning

- Random states
 - Bandits
 - ϵ -first
 - ϵ -greedy
 - Upper confidence bound
- Unobserved states
 - Robotics
 - Markov decision processes
 - Value / policy iteration
 - Linear programming implementation
 - Actor-Critic methods

This is a “for fun” lecture

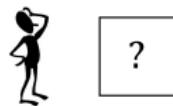
- Get some popcorn
- Put your feet up
- Enjoy!

The umbrella of online optimization

$\underset{x}{\text{minimize}} \quad g(x) \quad \text{where } g \text{ depends on something unknown}$

- **Unobserved state:** w is unknown

$\underset{x}{\text{minimize}} \quad f(x; w)$



- **Stochastic state:** $w \sim \mathcal{D}$ is a random variable

$\underset{x}{\text{minimize}} \quad \mathbb{E}_w[f(x; w)]$



- **Adversarial state:** w is adversarial

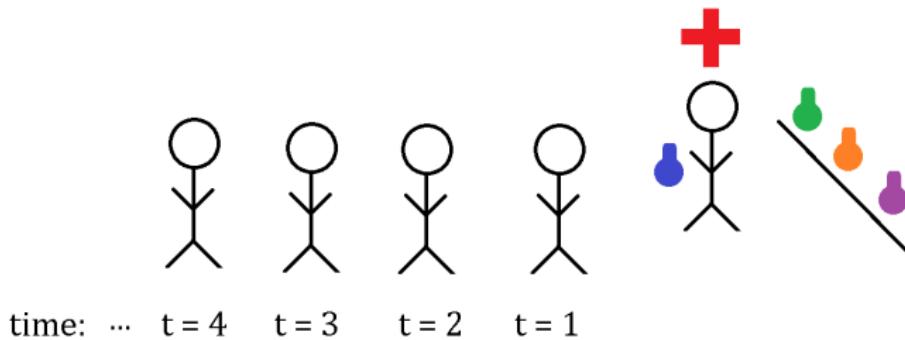
$\underset{x}{\text{minimize}} \quad \underset{w}{\max} \quad f(x; w)$



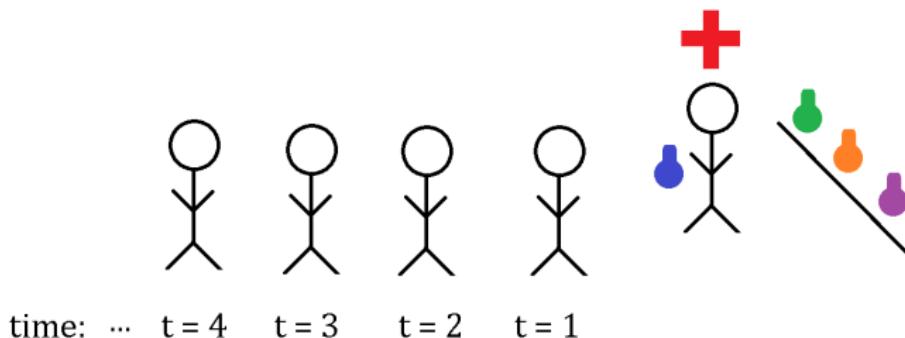
Stochastic online learning

A drug delivery problem

- m drugs, all super expensive, so only 1 drug per person
- Each person has a slightly different response to each drug
- Which drug should I recommend?



A drug delivery problem



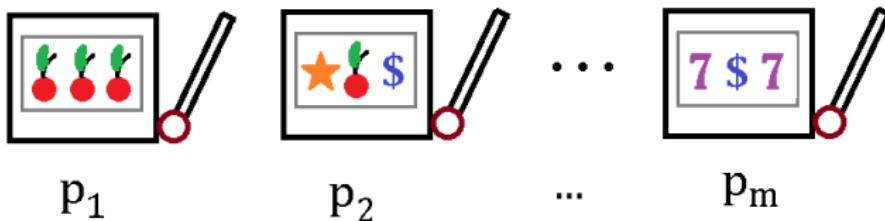
- Suppose each drug i has efficacy distributed as $\mathcal{N}(\mu_i, \sigma_i)$
- If I knew μ_i for $i = 1, \dots, m$ then I just give drug k where

$$k = \operatorname{argmax}_i \mu_i$$

- Best average efficacy

This is a (slightly harder) “bandit” problem

Bandit problem



- p_i = probability that arm i will win
- There is a ground truth p_i , but you do not know it
- The ground truth never changes
- General idea: estimate

$$p_i \approx \bar{p}_i = \frac{\# \text{ times arm } i \text{ wins}}{\# \text{ times I pull arm } i}$$

- Better guess if I pull it a lot, but that wastes money!

Regret

$$\text{Regret}(T) = \sum_{t=1}^T \text{loss}(A_t) - \text{loss}(A^*)$$

where

- A_t is the algorithm at time t
- A^* is the best algorithm possible

Sublinear regret:

$$\lim_{T \rightarrow \infty} \frac{\text{Regret}(T)}{T} \rightarrow 0$$

ϵ -first method

- Total T trials
- For ϵT trials, explore
 - Pull arm i with uniform probability
 - Update probability

$$\bar{p}_i = \frac{\# \text{ times arm } i \text{ wins}}{\# \text{ times I pull arm } i}$$



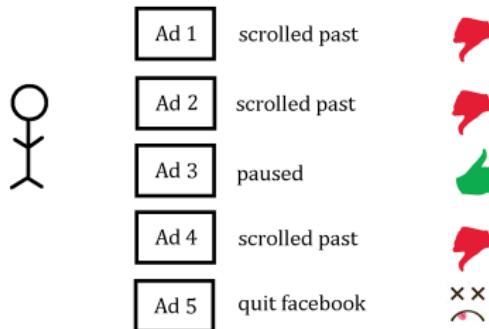
- For $(1 - \epsilon)T$ trials, exploit
 - Pull arm $i = \operatorname{argmax}_{i'} \bar{p}_{i'}$
- This also achieves sublinear regret!



ϵ -greedy method

- With probability ϵ , explore
 - Pull arm i with probability $1/m$ (uniform sampling)
- With probability $1 - \epsilon$, exploit
 - Pull arm $i = \operatorname{argmax}_{i'} \bar{p}_{i'}$
 - Break ties with uniform sampling
- Update probability
 - $$\bar{p}_i = \frac{\# \text{ times arm } i \text{ wins}}{\# \text{ times I pull arm } i}$$
- This simple strategy achieves sublinear regret!

Ad placement



- You show Joe an ad
- If Joe clicks or pauses on it → positive feedback!
- If Joe scrolls past it → negative feedback
- Limited budget: too many bad ads and Joe will stop paying attention

My ad trajectory

The New York Times
Sponsored

Fearlessness. Objectivity. Insight. All at this special rate.

THE NEW YORK TIMES SALE
Subscribe for just \$1 (Cdn) a week.
Ends soon.

The New York Times

WWW.NYT.COM
The facts. Up close.

Subscribe for \$1 (Cdn) a week.

151

79 Comments 6 Shares

Like Comment Share

My ad trajectory



Explore

My ad trajectory

prime video Amazon Prime Video ***

Sponsored · 45

What will Hanna find beyond the forest? Find out on March 29, only on PrimeVideo.com.

HANNA prime video



I want to know what's beyond the forest.

PRIMEVIDEO.COM Watch More

Start your free trial now

45 7 Comments 7 Shares 868K Views

Like Comment Share 14

My ad trajectory

 **NPR Music**
32 mins · ⓘ

A New York State Supreme Court judge issued the decision almost exactly a year after the famed opera house fired the former music director James Levine.



NPR.ORG

Majority Of James Levine's Defamation Claims Against Met Opera Dismissed

2

Like Comment Share

Explore

My ad trajectory



Exploit?

My ad trajectory



Exploit?

My ad trajectory

Crest
Sponsored

Introducing new Crest Whitening Therapy Toothpaste with Charcoal. Gentle Whitening. Premium Ingredients.

The image shows a Facebook post from the brand 'Crest'. The post is sponsored and features a product image of a tube of Crest Whitening Therapy toothpaste with charcoal. The text 'GENTLE WHITENING.' is displayed above the product. A play button icon is overlaid on the product image, indicating it's a video advertisement. Below the image, the text 'AMAZON.CA' and a 'Shop Now' button are visible. At the bottom of the post, there are engagement metrics: 294 likes, 62 comments, and 49 shares. There are also standard social media interaction buttons for 'Like', 'Comment', and 'Share'.

GENTLE WHITENING.

AMAZON.CA

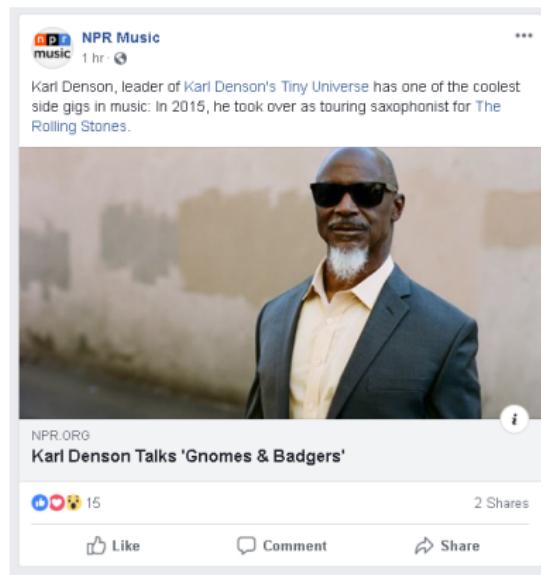
Shop Now

294

62 Comments 49 Shares

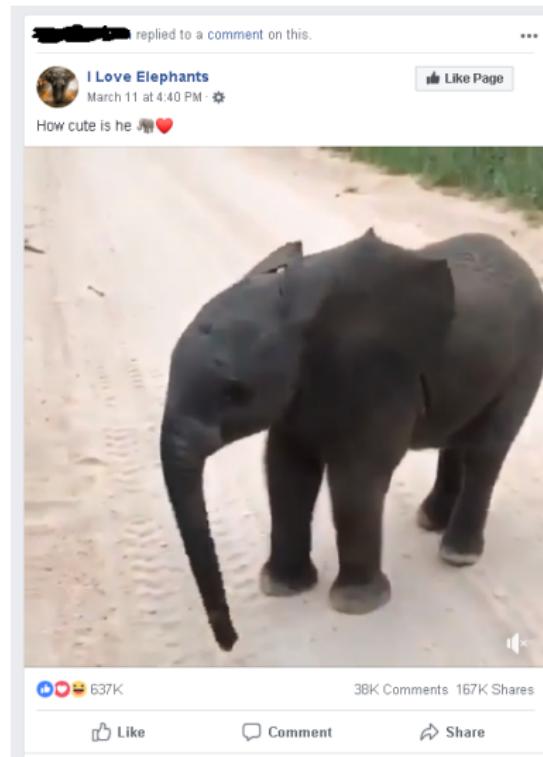
Like Comment Share

My ad trajectory



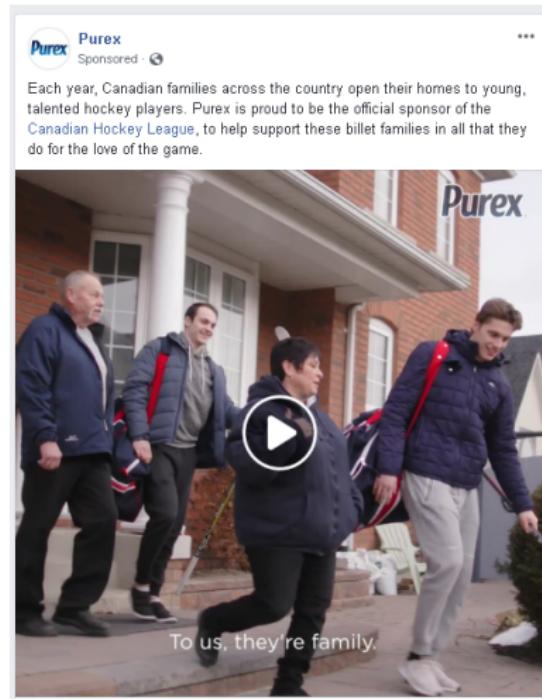
Exploit?

My ad trajectory



Explore (contextual?)
20

My ad trajectory



568

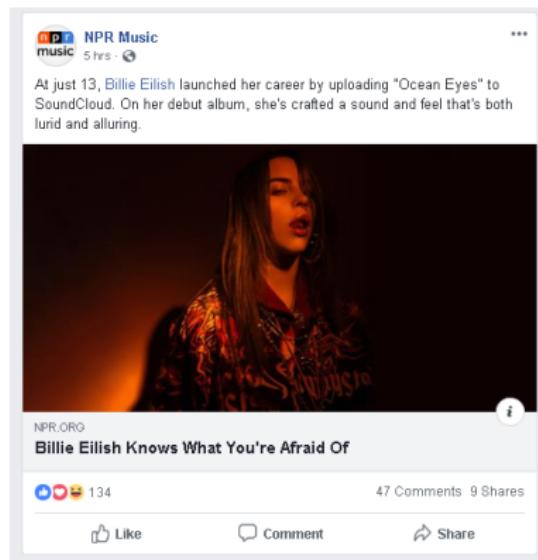
40 Comments 84 Shares 740K Views

Like

Comment

Share

My ad trajectory



Exploit

My ad trajectory



Exploit

My ad trajectory

 Seattle Opera
Sponsored · 3

In honor of the children's chorus for our upcoming opera Carmen, which includes 16 youth singers from the Pacific Northwest, we're celebrating a history of youth performers at Seattle Opera. Our grand new production of Carmen features lavish scenery, traditional costumes, and some of the most famous classical music ever written.

Carmen plays May 4-19, 2019 at McCaw Hall.
Tickets & info: seattleopera.org/carmen



Payton Barnes - Porgy and Bess
Learn More
Photo: Philip Newton

Chloe Carrasquillo - The Magic Flute
Learn More
Photo: Genevieve Hathaway

Like Comment Share

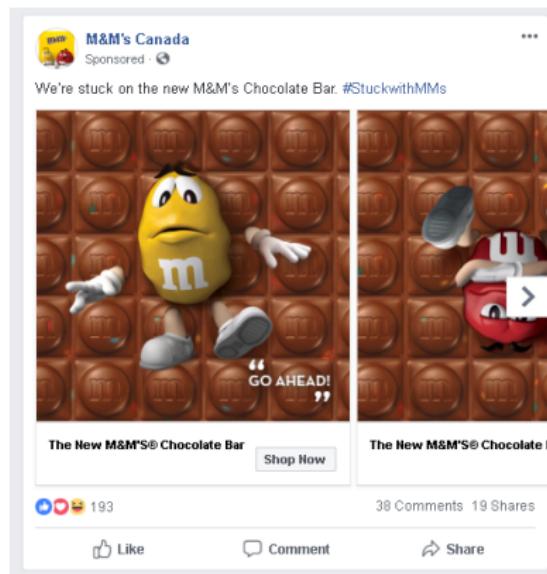
Explore

My ad trajectory



Exploit

My ad trajectory



Explore

My ad trajectory



Exploit

Upper confidence bound (UCB)

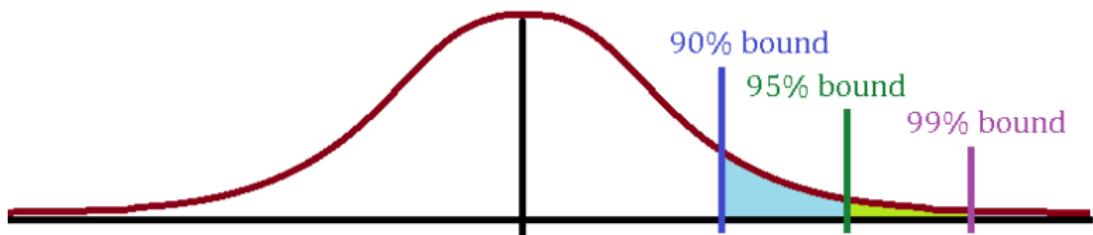
- Consider X a random variable with mean μ and variance σ^2
- After k observations x_1, \dots, x_k , the sample mean

$$\bar{x} = \frac{1}{k} \sum_{i=1}^k x_i$$

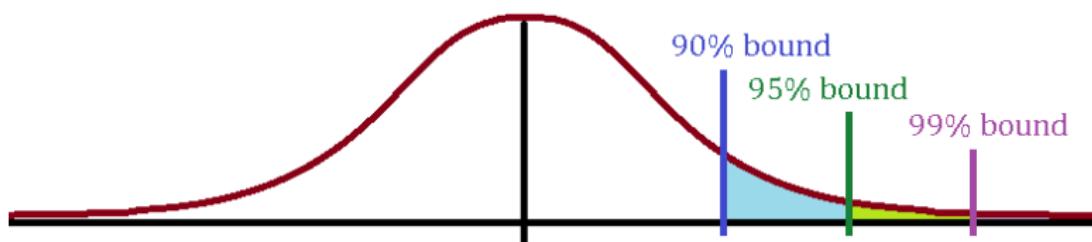
is also a random variable, with mean \bar{x} and variance σ^2/k

- The upper confidence bound for $1-\alpha$ confidence is d where

$$\Pr(\bar{x} > d) \leq \alpha$$

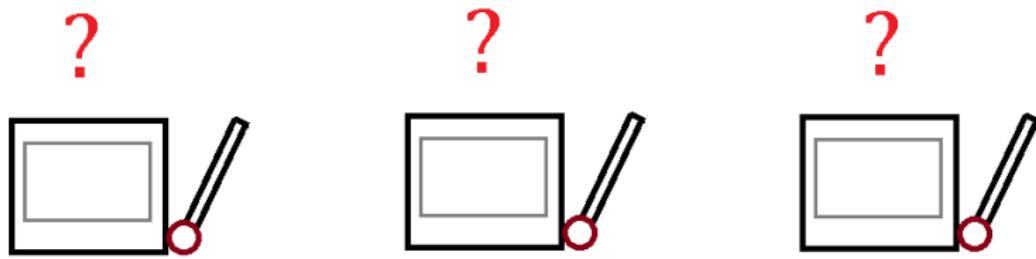


Upper confidence bound method



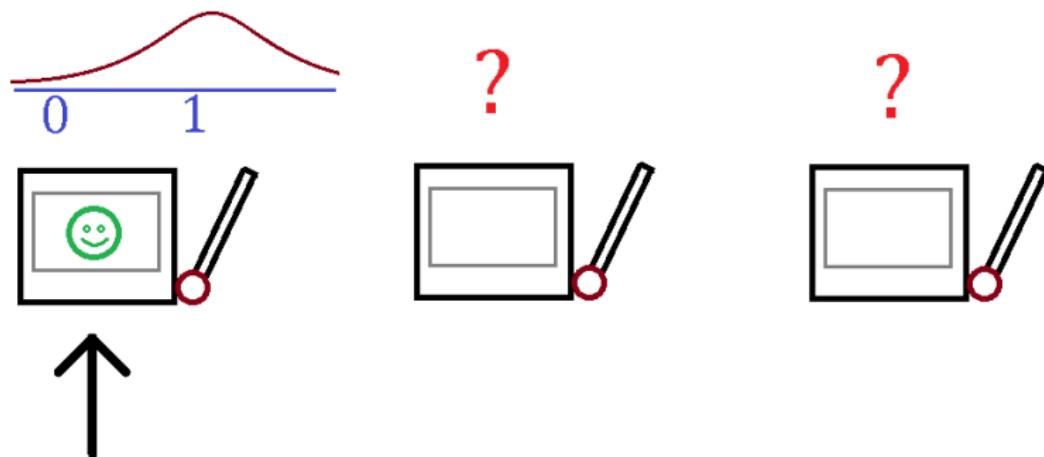
- If you don't know anything, pull any arm
- Now pick the arm with the largest UCB for some confidence $1 - \alpha_t$
- α_t decays to 0
 - never stops exploring, but happens more seldomly for not-best arm
- an “optimistic” method

Upper confidence bound method



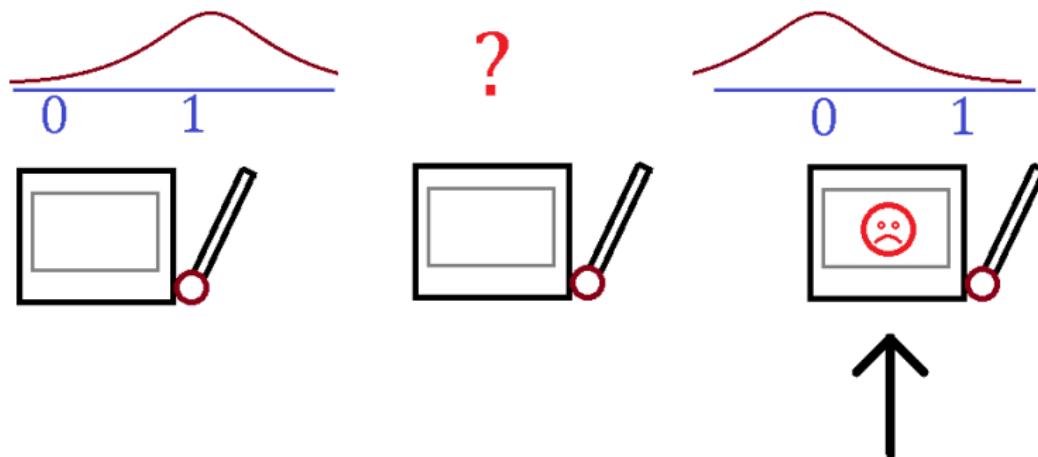
Regret = total pulls - total wins

Upper confidence bound method



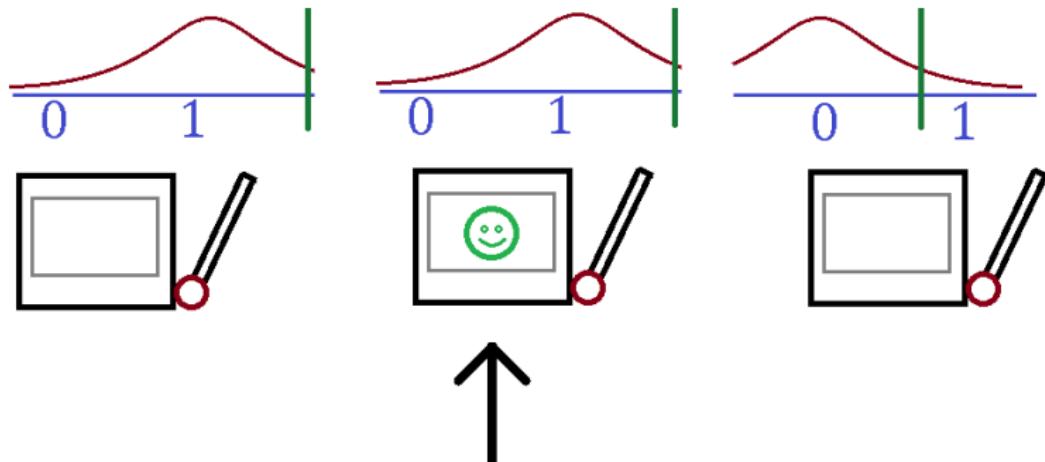
$$\text{Regret} = 1 - 1$$

Upper confidence bound method



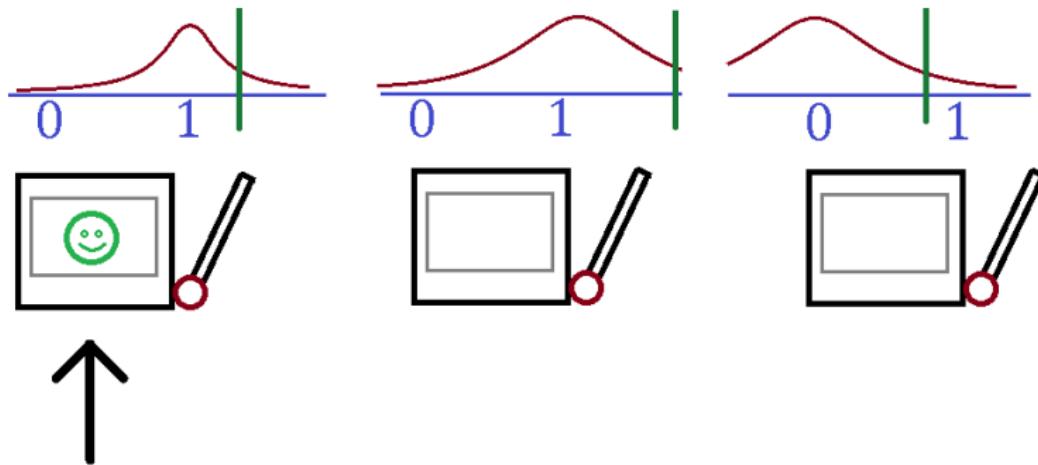
$$\text{Regret} = 2 - 1$$

Upper confidence bound method



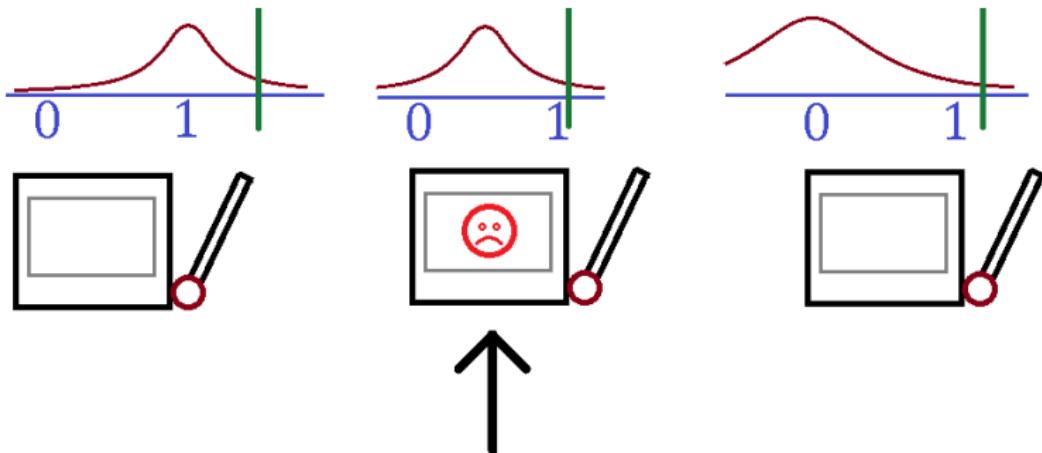
$$\text{Regret} = 3 - 1$$

Upper confidence bound method



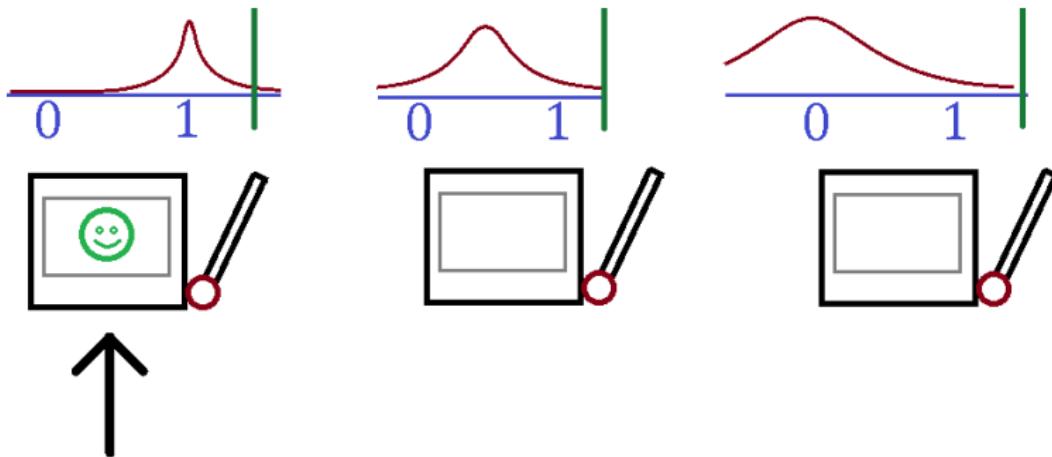
$$\text{Regret} = 4 - 2$$

Upper confidence bound method



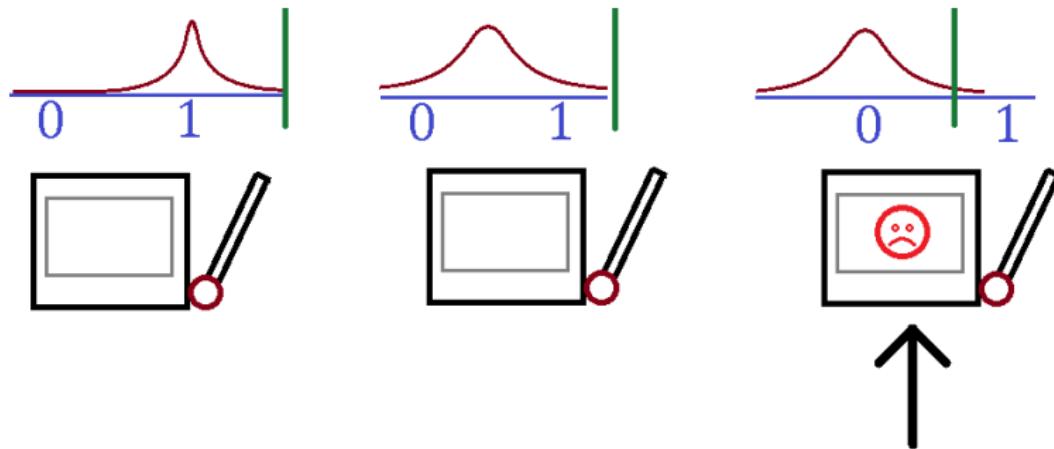
$$\text{Regret} = 5 - 3$$

Upper confidence bound method



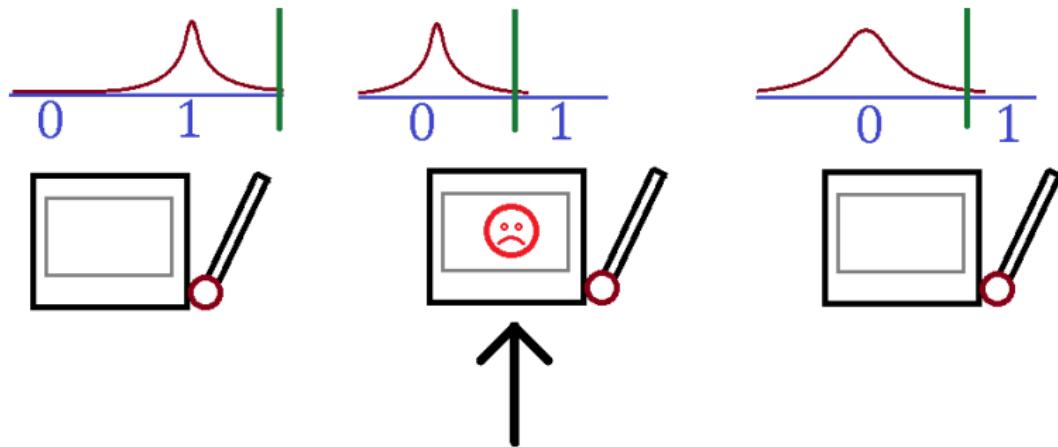
$$\text{Regret} = 6 - 4$$

Upper confidence bound method



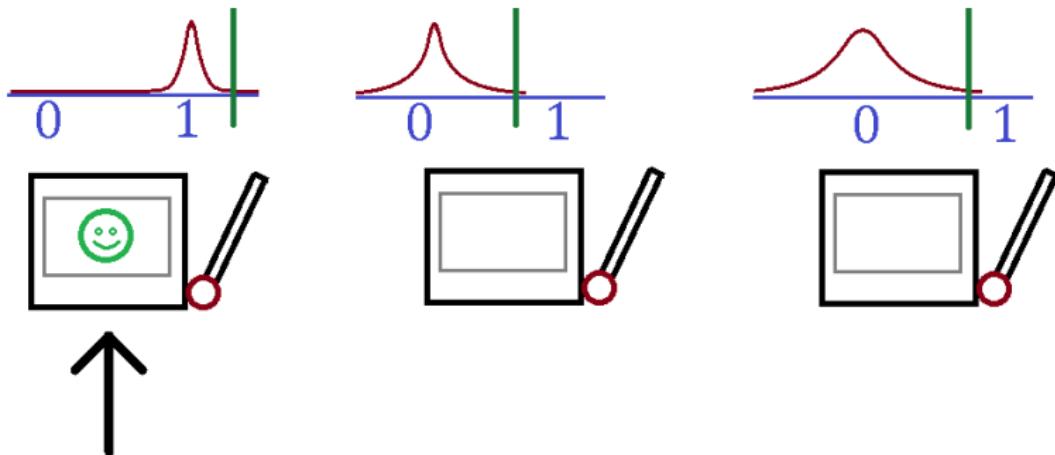
$$\text{Regret} = 7 - 4$$

Upper confidence bound method



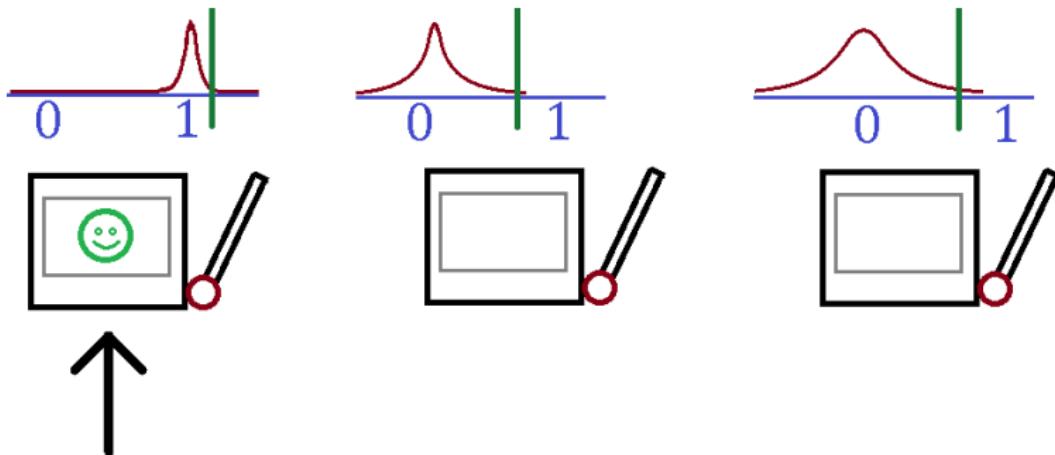
$$\text{Regret} = 8 - 4$$

Upper confidence bound method



$$\text{Regret} = 9 - 5$$

Upper confidence bound method



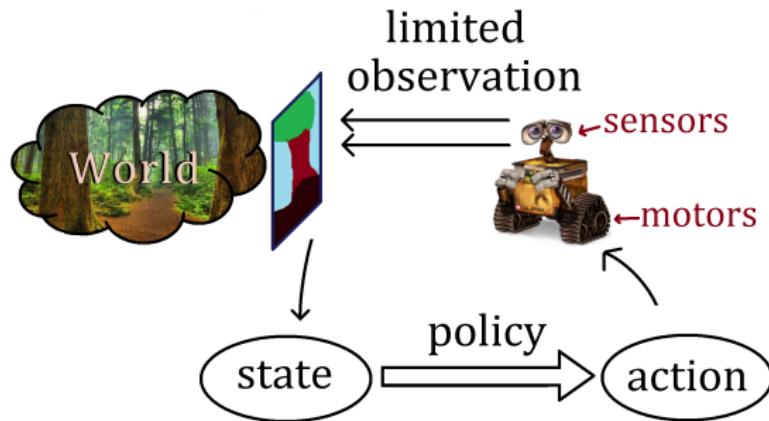
$$\text{Regret} = 10 - 5$$

Unobserved states

Spot the dog



Control cycle

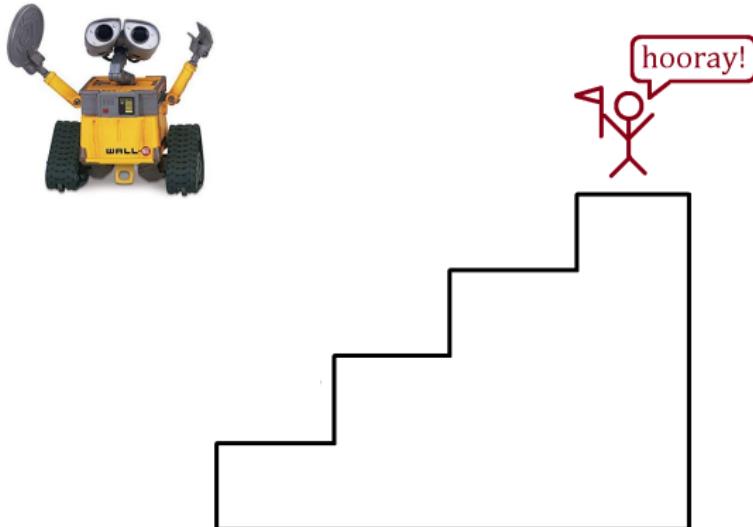


- Robot glimpses world → state updated
- Policy → Robot moves
- Robot gets a new glimpse → state updated
- Rinse and repeat

Robot climbing stairs

POLICY

Move to state of higher reward



What's a good policy?

Goal: Robot climbing the stairs

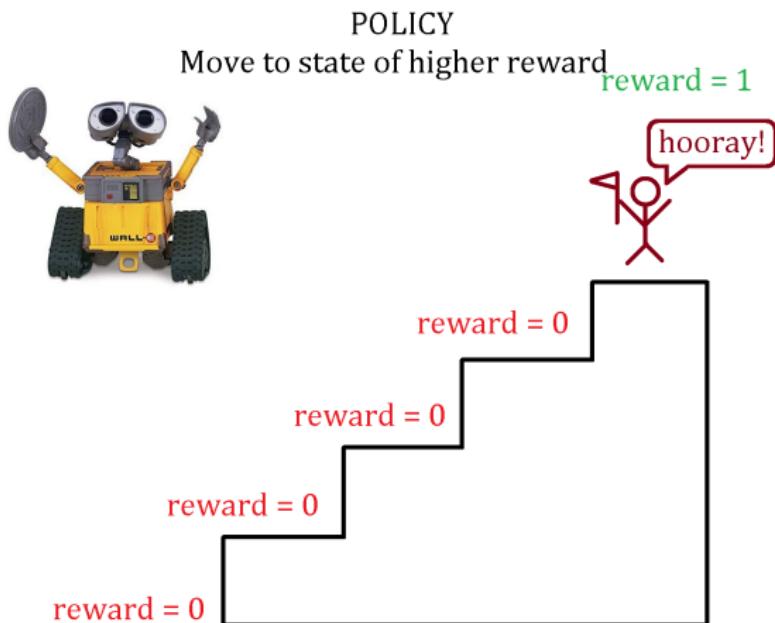
- Markov state transition:(past given current doesn't matter)

$\Pr(s^{(t+1)} \mid s^{(t)}, a^{(t)})$ fixed but unknown

- Reward:

$$r(s) = \begin{cases} 1 & \text{if } s = \text{top step} \\ 0 & \text{otherwise} \end{cases}$$

What's a good policy?



What's a good policy?

Goal: Robot climbing the stairs

- Markov state transition:(past given current doesn't matter)

$$\Pr(s^{(t+1)} \mid s^{(t)}, a^{(t)}) \text{ fixed but unknown}$$

- Reward:

$$r(s) = \begin{cases} 1 & \text{if } s = \text{top step} \\ 0 & \text{otherwise} \end{cases}$$

- Policy (first try)

$$\pi(s) = a = \operatorname{argmax}_a \{\mathbb{E}_{s' \mid s, a}[r(s')]\}$$

- What's wrong with this picture?

What's a good policy?

Goal: Robot climbing the stairs

- Markov state transition:(past given current doesn't matter)

$$\Pr(s^{(t+1)} \mid s^{(t)}, a^{(t)}) \text{ fixed but unknown}$$

- Reward:

$$r(s) = \begin{cases} 1 & \text{if } s = \text{top step} \\ 0 & \text{otherwise} \end{cases}$$

- Policy (first try)

$$\pi(s) = a = \operatorname{argmax}_a \{\mathbb{E}_{s' \mid s, a}[r(s')]\}$$

- What's wrong with this picture? Ans: reward is too sparse!

What's a good policy?

Goal: Robot climbing the stairs

- Markov state transition:

$$s^{(t+1)} = \Pr(s^{(t+1)} \mid s^{(t)}, a^{(t)})$$

- Reward:

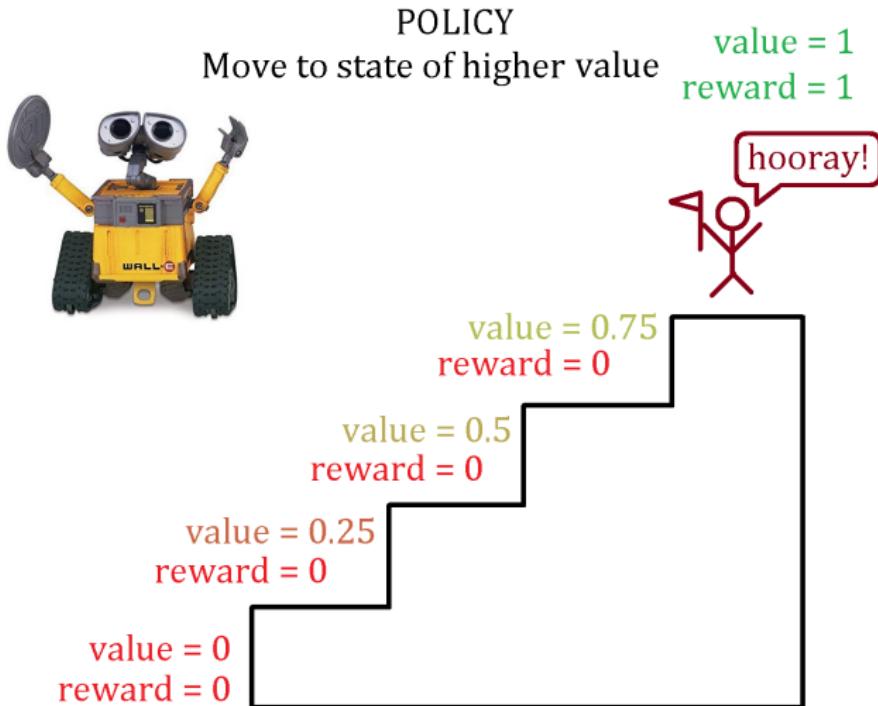
$$r(s) = \begin{cases} 1 & \text{if } s = \text{top step} \\ 0 & \text{otherwise} \end{cases}$$

- Value given policy: (discounted reward)

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^t r(s^{(t)}) : s^{(0)} = s, s^{(t)} = \pi(s^{(t-1)}) \right], \quad 0 < \gamma < 1$$

- Interpretation: How well will I do in the near future, if I follow this policy?

What's a good policy?



What's a good policy?

- Value given policy: (discounted reward)

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^t r(s^{(t)}) : s^{(0)} = s, s^{(t)} = \pi(s^{(t-1)}) \right], \quad 0 < \gamma < 1$$

- Optimal policy

$$\pi^*(s) = \operatorname{argmax}_{\pi} V^\pi(s)$$

- “Best” value function

$$V^*(s) = V^{\pi^*}$$

- Chicken and egg problem?

Bellman condition

- Bellman optimality equation: $V^*(s)$ uniquely solves

$$V^*(s) = r(s) + \gamma \max_a \sum_{s'} \Pr(s' | s, a) V^*(s')$$

- Optimal policy

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} \Pr(s' | s, a) V^*(s')$$

Computing optimal policy and value function

Value iteration

$$\begin{aligned} V^{(0)}(s) &= 0 \quad \forall s \\ V^{(t+1)}(s) &= R(s) + \gamma \max_a \sum_{s'} \Pr(s' | s, a) \quad \forall s \end{aligned}$$

Then $V^{(t)} \xrightarrow{t \rightarrow \infty} V^*$

Policy iteration

$$\begin{aligned} \pi^{(0)}(s) &= 0 \quad \forall s \\ V^{(t)}(s) &= R(s) + \gamma \sum_{s'} \Pr(s' | s, \pi(s)) V^{(t)}(s) \text{ through roll-out} \\ \pi^{(t+1)}(s) &= \operatorname{argmax}_a \sum_{s'} \Pr(s' | s, a) \quad \forall s \end{aligned}$$

Then $\pi^{(t)} \xrightarrow{t \rightarrow \infty} \pi^*$

Linear programming

V^* is the solution to the following infinite-dimensional linear program

$$(P) \quad \begin{aligned} & \text{minimize}_V && \sum_s V(s) \\ & \text{subject to} && V(s) \geq R(s) + \gamma \sum_{s'} \Pr(s' | s, a) V(s') \quad \forall a \end{aligned}$$

Why does this work?

- Feasibility implies

$$V(s) \geq R(s) + \gamma \max_a \sum_{s'} \Pr(s' | s, a) V(s')$$

- Optimality implies tightness

$$V(s) = R(s) + \gamma \max_a \sum_{s'} \Pr(s' | s, a) V(s')$$

- This means Bellman's condition is satisfied

Dual linear programming

V^* is the solution to the following infinite-dimensional linear program

$$(D) \quad \begin{aligned} & \underset{\mu}{\text{maximize}} && \sum_s \sum_a R(s) \mu(s, a) \\ & \text{subject to} && \sum_{a'} \mu(s', a) = e + \gamma \sum_s \sum_a \Pr(s' | s, a) \mu(s, a) \quad \forall s' \\ & && \mu(s, a) \geq 0 \quad \forall s, a \end{aligned}$$

Interpretation

- μ is discounted probability

$$\mu(s, a) = \sum_{t=0}^{\infty} \gamma^t \Pr(s^{(t)} = s, a^{(t)} = a)$$

- Optimal policy

$$\pi^*(s) = \max_a \mu(s, a)$$

Caveats

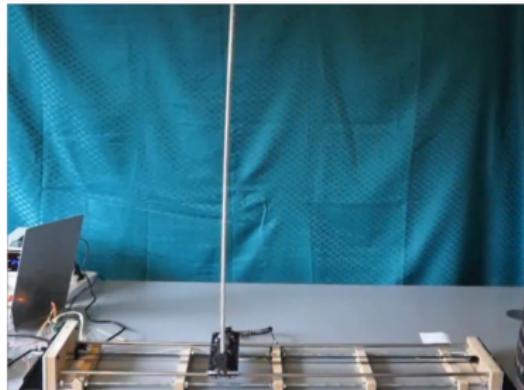
- In general, you cannot really optimize over infinite spaces
 - Discretize state space into grid
 - Discretize functions with basis functions

$$V(s) = \sum_{i=1}^n \alpha_i b_i(s)$$

b_i is fixed, learn α_i

- Even if your state space was finite, a robot can only see its current state
 - Q-learning, SARSA, ...
 - At each iteration, $V(s)$ is only updated for one value of s
 - Alternatively, roll-out (check all s by playing game many times)

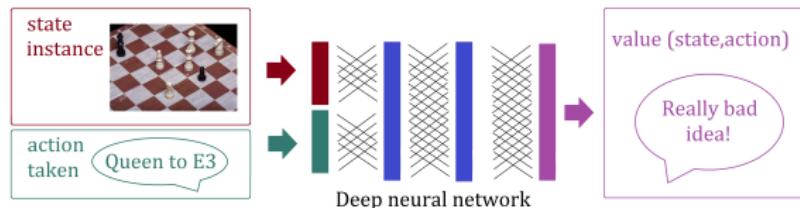
Cart Pole using temporal difference learning



<https://www.youtube.com/watch?v=5Q14EjnOJZc>

Deep reinforcement learning: Actor-critic

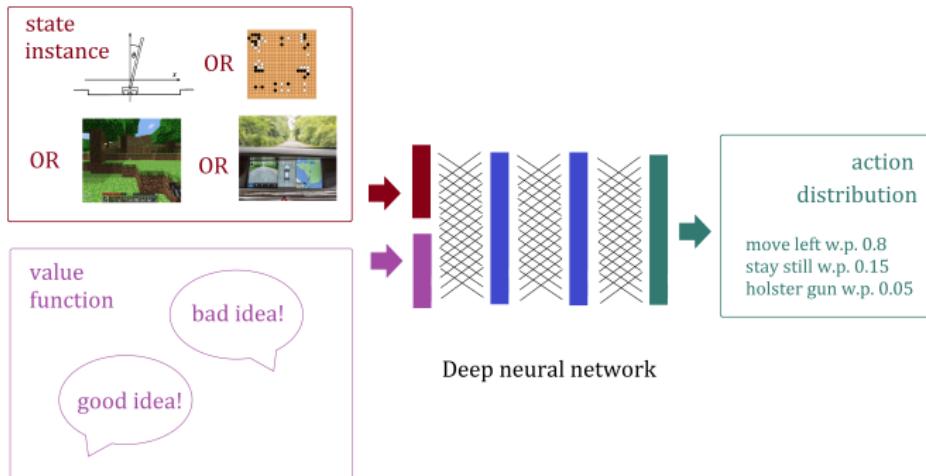
Critic



- Determines the value function given state and action
 - Also called Q function
- Function is approximated by a deep learning network
- This is a good idea when states are represented by images

Deep reinforcement learning: Actor-critic

Actor



- Determines action probability distribution, given state and value function
- Successes: AlphaGo, Atari, Doom, ...

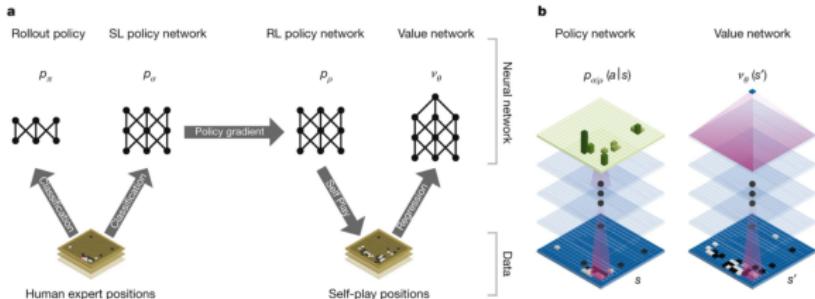
Google's AlphaGo AI wins three-match series against the world's best Go player



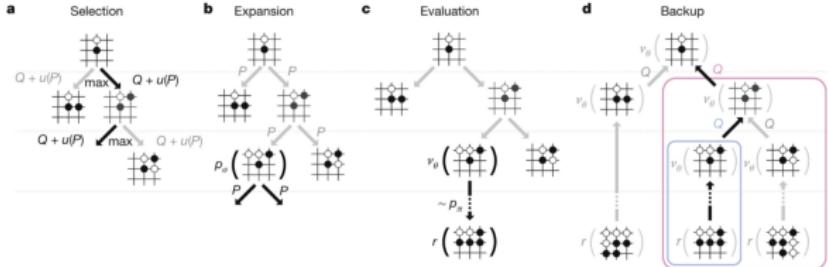
Google's AlphaGo AI has once again made the case that machines are now smarter than man — when it comes to games of strategy, at least.



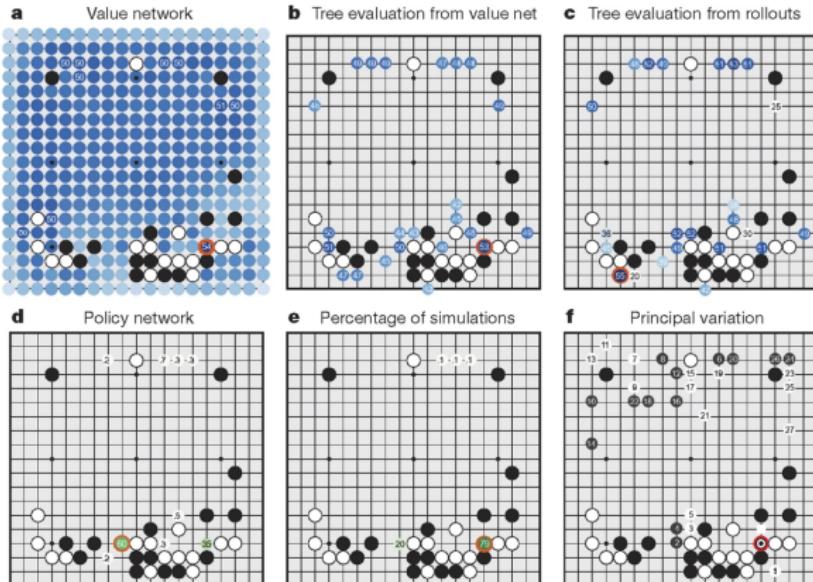
All games of perfect information have an optimal value function, $v^*(s)$, which determines the outcome of the game, from every board position or state s , under perfect play by all players. These games may be solved by recursively computing the optimal value function in a search tree containing approximately b^d possible sequences of moves, where b is the game's breadth (number of legal moves per position) and d is its depth (game length). In large games, such as chess ($b \approx 35, d \approx 80$)¹ and especially Go ($b \approx 250, d \approx 150$)¹, exhaustive search is infeasible^{2,3}, but the effective search space can be reduced by two general principles. First, the depth of the search may be reduced by position evaluation: truncating the search tree at state s and replacing the subtree below s by an approximate value function $v(s) \approx v^*(s)$ that predicts the outcome from state s . This approach has led to superhuman performance in chess⁴, checkers⁵ and othello⁶, but it was believed to be intractable in Go due to the complexity of the game⁷. Second, the breadth of the search may be reduced by sampling actions from a policy $p(a|s)$ that is a probability distribution over possible moves a in position s . For example, Monte Carlo rollouts⁸ search to maximum depth without branching at all, by sampling long sequences of actions for both players from a policy p . Averaging over such rollouts can provide an effective position evaluation, achieving superhuman performance in backgammon⁸ and Scrabble⁹, and weak amateur level play in Go¹⁰.



a, A fast rollout policy p_π and supervised learning (SL) policy network p_σ are trained to predict human expert moves in a data set of positions. A reinforcement learning (RL) policy network p_ρ is initialized to the SL policy network, and is then improved by policy gradient learning to maximize the outcome (that is, winning more games) against previous versions of the policy network. A new data set is generated by playing games of self-play with the RL policy network. Finally, a value network v_θ is trained by regression to predict the expected outcome (that is, whether the current player wins) in positions from the self-play data set. **b**, Schematic representation of the neural network architecture used in AlphaGo. The policy network takes a representation of the board position s as its input, passes it through many convolutional layers with parameters σ (SL policy network) or ρ (RL policy network), and outputs a probability distribution $p_\sigma(a|s)$ or $p_\rho(a|s)$ over legal moves a , represented by a probability map over the board. The value network similarly uses many convolutional layers with parameters θ , but outputs a scalar value $v_\theta(s')$ that predicts the expected outcome in position s' .



a. Each simulation traverses the tree by selecting the edge with maximum action value Q , plus a bonus $u(P)$ that depends on a stored prior probability P for that edge. **b.** The leaf node may be expanded; the new node is processed once by the policy network p_a and the output probabilities are stored as prior probabilities P for each action. **c.** At the end of a simulation, the leaf node is evaluated in two ways: using the value network v_θ ; and by running a rollout to the end of the game with the fast rollout policy p_n , then computing the winner with function r . **d.** Action values Q are updated to track the mean value of all evaluations $r(\cdot)$ and $v_\theta(\cdot)$ in the subtree below that action.



For each of the following statistics, the location of the maximum value is indicated by an orange circle. **a**, Evaluation of all successors s' of the root position s , using the value network $v_\theta(s')$; estimated winning percentages are shown for the top evaluations. **b**, Action values $Q(s, a)$ for each edge (s, a) in the tree from root position s ; averaged over value network evaluations only ($\lambda = 0$). **c**, Action values $Q(s, a)$, averaged over rollout evaluations only ($\lambda = 1$). **d**, Move probabilities directly from the SL policy network, $p_\theta(a|s)$; reported as a percentage (if above 0.1%). **e**, Percentage frequency with which actions were selected from the root during simulations. **f**, The principal variation (path with maximum visit count) from AlphaGo's search tree. The moves are presented in a numbered sequence. AlphaGo selected the move indicated by the red circle; Fan Hui responded with the move indicated by the white square; in his post-game commentary he preferred the move (labelled 1) predicted by AlphaGo.

Configuration and strength^[60]

Versions	Hardware	Elo rating	Date	Results
AlphaGo Fan	176 GPUs, ^[51] distributed	3,144 ^[50]	Oct 2015	5:0 against Fan Hui
AlphaGo Lee	48 TPUs, ^[51] distributed	3,739 ^[50]	Mar 2016	4:1 against Lee Sedol
AlphaGo Master	4 TPUs, ^[51] single machine	4,858 ^[50]	May 2017	60:0 against professional players; Future of Go Summit
AlphaGo Zero (40 block)	4 TPUs, ^[51] single machine	5,185 ^[50]	Oct 2017	100:0 against AlphaGo Lee 89:11 against AlphaGo Master
AlphaZero (20 block)	4 TPUs, single machine	5,018 ^[61]	Dec 2017	60:40 against AlphaGo Zero (20 block)

<https://en.wikipedia.org/wiki/AlphaGo>

Computer plays doom



<https://www.youtube.com/watch?v=oo0TraGu6QY>

Computer plays doom

YouTube computer plays doom

4 Played by Deviant's Choplif

Matthew Laro 2 years ago

Are the other players AI controlled using raw pixel data or in-game bots?

dr 22 41 REPLY

Hide replies ▾

• Deviant's Choplif 2 years ago

Others are in game bots which use internal game information.

dr 10 41 REPLY

• n 2 years ago

At 1:40 it was going to shoot with shotgun when it noticed he killed them all, interesting...

dr 24 41 REPLY

Hide replies ▾

• Pernarox 2 years ago

That was likely because the death animation was hidden from it through the player in front

dr 10 41 REPLY

• Metal Aromatic 2 years ago

You be careful out there, and that will create Skynet

dr 41 41 REPLY

View 3 replies ▾

• Net 8 M 2 years ago

Standing still while shooting would not work against real players.

dr 31 41 REPLY

Hide replies ▾

• robot stamp 2 years ago

exactly what I thought - no movement when shooting, no strafing, no backwards. From the article I imagined the AI would be tougher, the bots wondering the map are also very stupid.

dr 2 41 REPLY

• Random Dumb 2 years ago

It seems like it's using a simplified control scheme, possibly to minimize garbage outputs. It doesn't seem like the AI was given awareness of the strafe keys - it probably is only aware of forward, left, right, and shoot, from what I can see in the video.

<https://www.youtube.com/watch?v=oo0TraGu6QY>

Some more references

- Online convex optimization and no-regret learning: algorithms, guarantees, and applications by Belmega et al.
- Blogs:
 - First part: google stochastic bandits, ϵ -greedy, sublinear regret learning
 - Second part: google Reinforcement learning, Q-learning, actor-critic, deep reinforcement learning
- LP formulation of value / policy optimization:
 - Class slides: <http://www.cs.cmu.edu/afs/cs/academic/class/15780-s16/www/slides/mdps.pdf>