

Strategic Interactive Decision-Making

Thomas Kleine Buening

The Alan Turing Institute & University of Oxford

December 5, 2024

**The
Alan Turing
Institute**



UNIVERSITY OF
OXFORD



INTENTION
REDUCE COBRA
POPULATION



INTENTION
REDUCE COBRA
POPULATION



ACTION
A BOUNTY FOR
DEAD COBRAS!



INTENTION

REDUCE COBRA
POPULATION



ACTION

A BOUNTY FOR
DEAD COBRAS!



EFFECT

PEOPLE START
COBRA FARMING



INTENTION
REDUCE COBRA
POPULATION



ACTION
A BOUNTY FOR
DEAD COBRAS!



EFFECT
PEOPLE START
COBRA FARMING

***“Anything that can go wrong
will go wrong.”***

Murphy's Law



INTENTION
REDUCE COBRA
POPULATION



ACTION
A BOUNTY FOR
DEAD COBRAS!



EFFECT
PEOPLE START
COBRA FARMING

***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur

***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



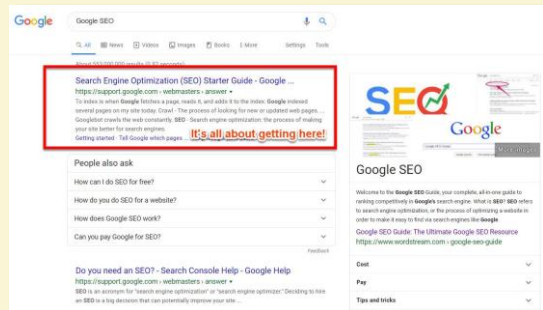
***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



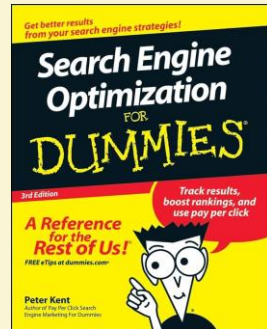
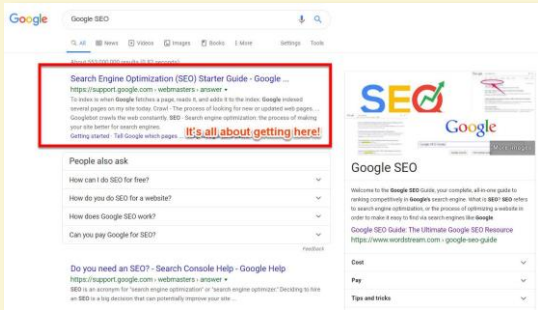
***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur



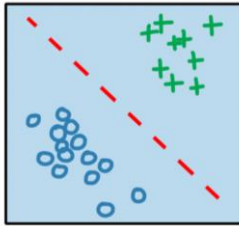
***“Any system that can be gamed
will be gamed.”***

W. Brian Arthur

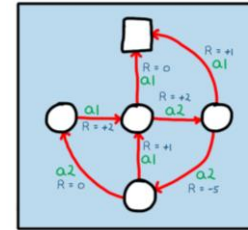
We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

Supervised Learning

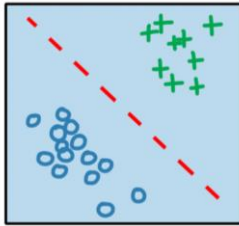


Reinforcement Learning



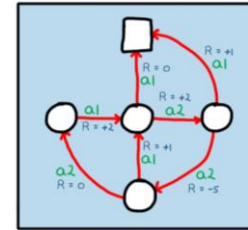
We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

Supervised Learning



Adversarial Robustness

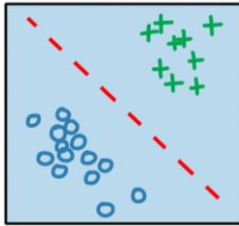
Reinforcement Learning



Adversarial Robustness

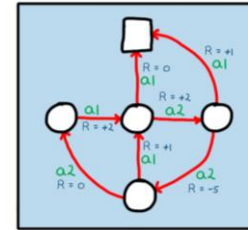
We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

Supervised Learning



Adversarial Robustness

Reinforcement Learning

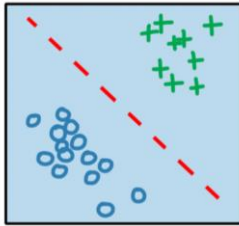


Adversarial Robustness



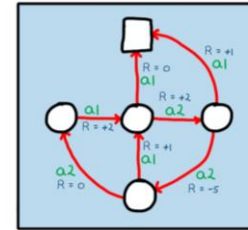
We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

Supervised Learning



Adversarial Robustness
Strategic Classification

Reinforcement Learning

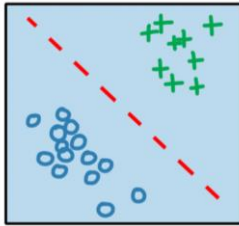


Adversarial Robustness
Corruption-Robust RL



We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

Supervised Learning

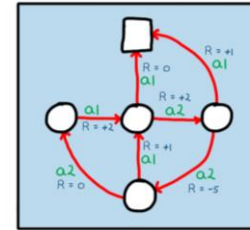


Adversarial Robustness
Strategic Classification



(1) (2)

Reinforcement Learning

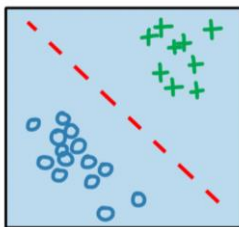


Adversarial Robustness
Corruption-Robust RL

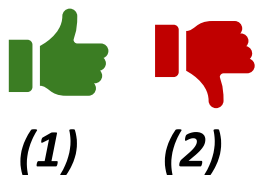


We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

Supervised Learning

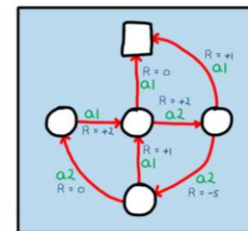


Adversarial Robustness
Strategic Classification



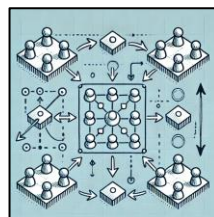
(1) (2)

Reinforcement Learning



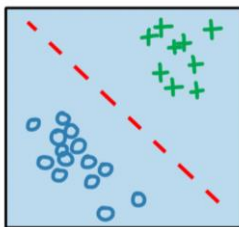
Adversarial Robustness
Corruption-Robust RL

Mechanism Design

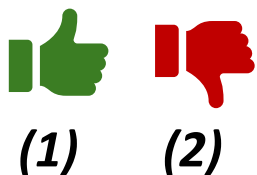


We want to design **ML algorithms** that are (1) **robust** against **strategic behavior** and (2) **incentivize agent behavior** that is aligned with the **system's goals**.

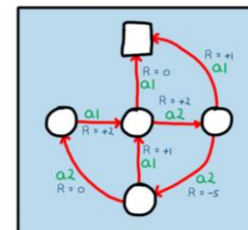
Supervised Learning



Adversarial Robustness
Strategic Classification



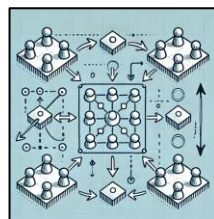
Reinforcement Learning



Adversarial Robustness
Corruption-Robust RL

Strategic Interactive Decision-Making

Mechanism Design



Strategic Linear Contextual Bandits

joint work with Aadirupa Saha, Christos Dimitrakakis, Haifeng Xu

*recommends the
channels' content*

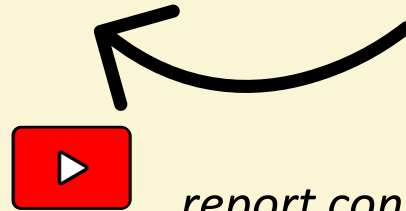


maximizing individual exposure / profit



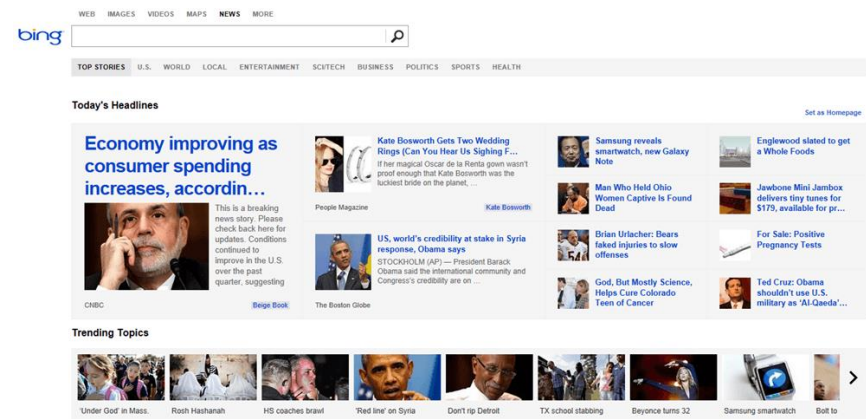
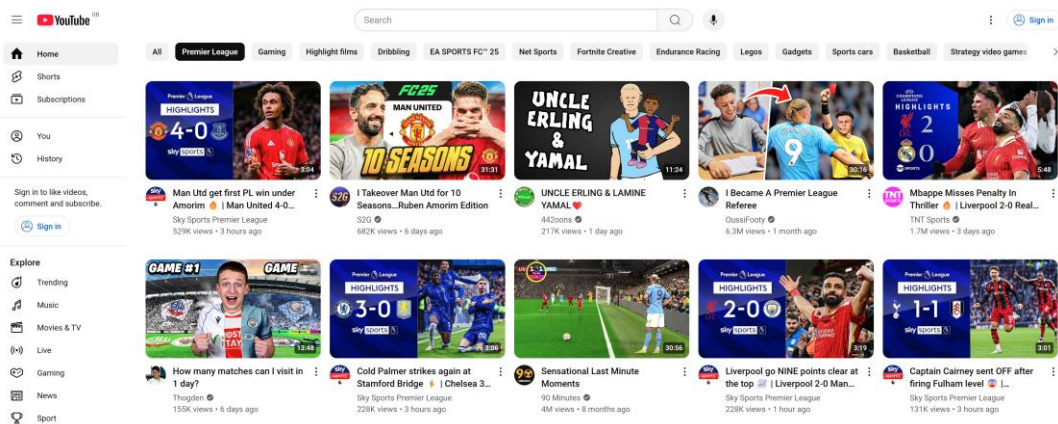
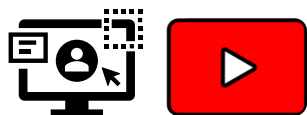
maximizing platform performance

- (1) make good recommendations
- (2) incentivize good content / truthfulness



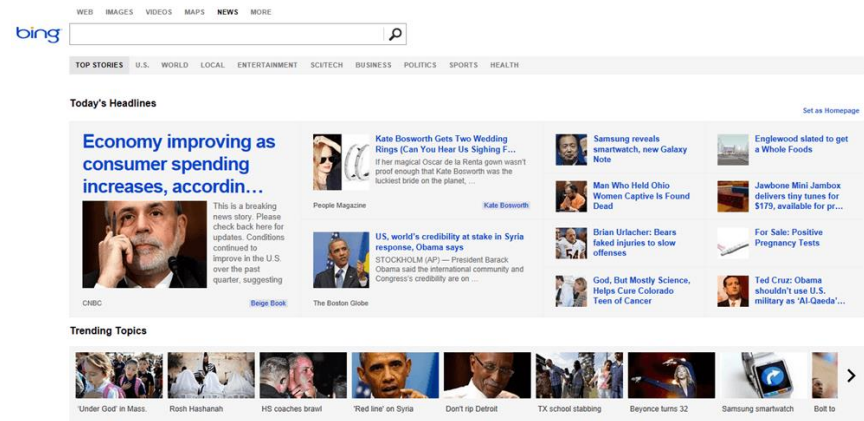
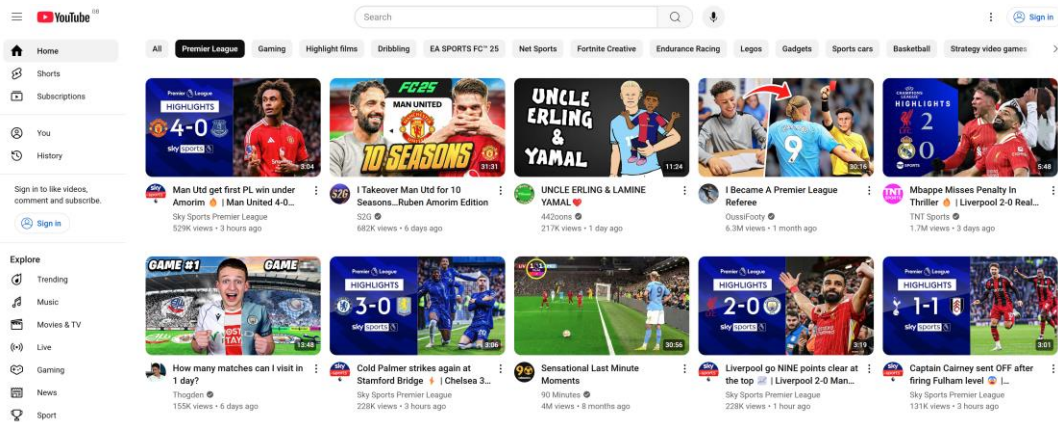
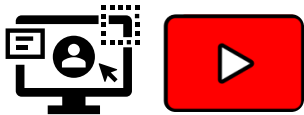
*report content
and data strategically*

Linear Contextual Bandits




Linear Contextual Bandits

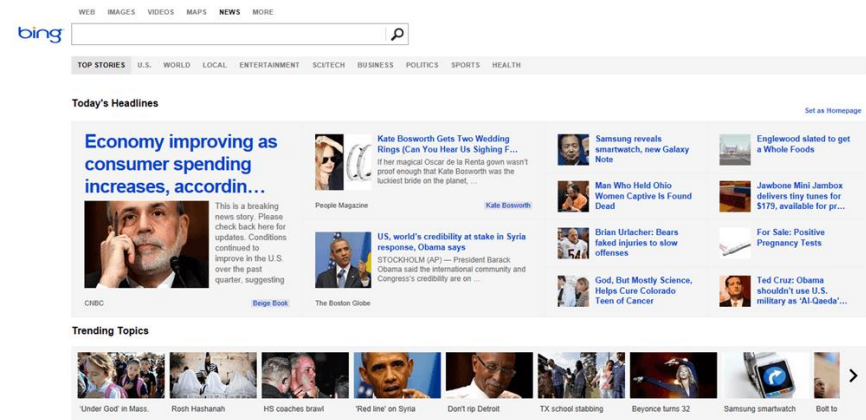
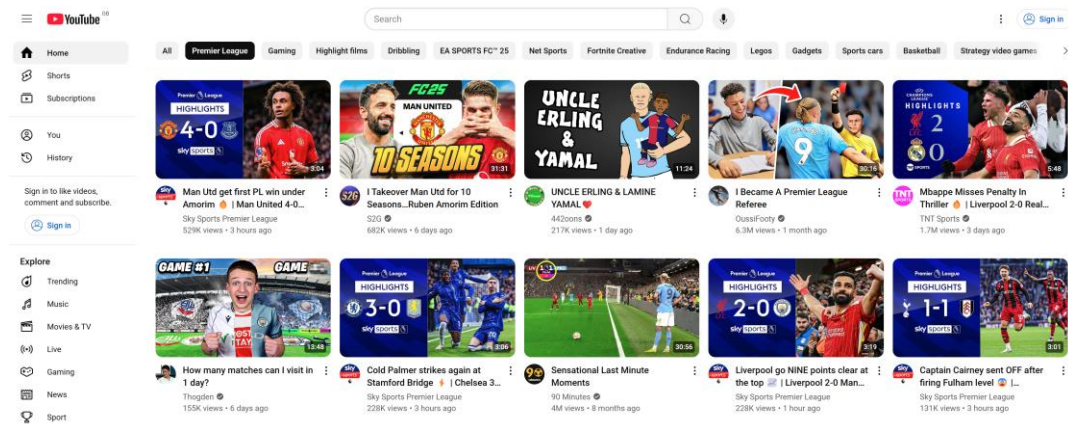
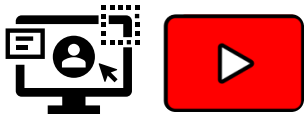
- T rounds, K arms ()



Linear Contextual Bandits

- T rounds, K arms ()

For $t = 1, \dots, T$:

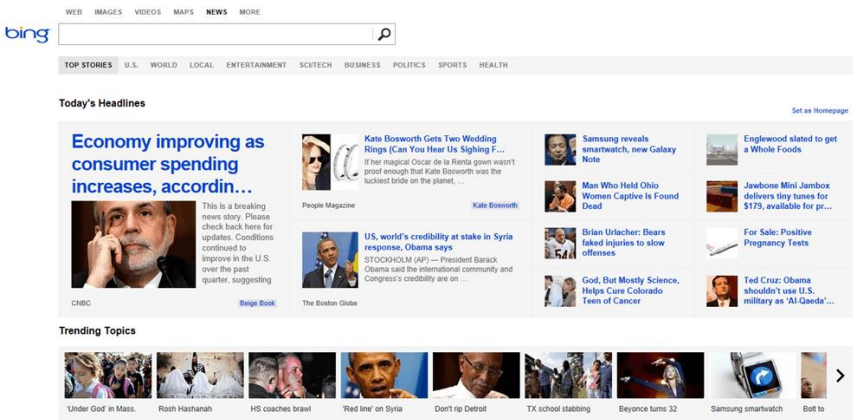
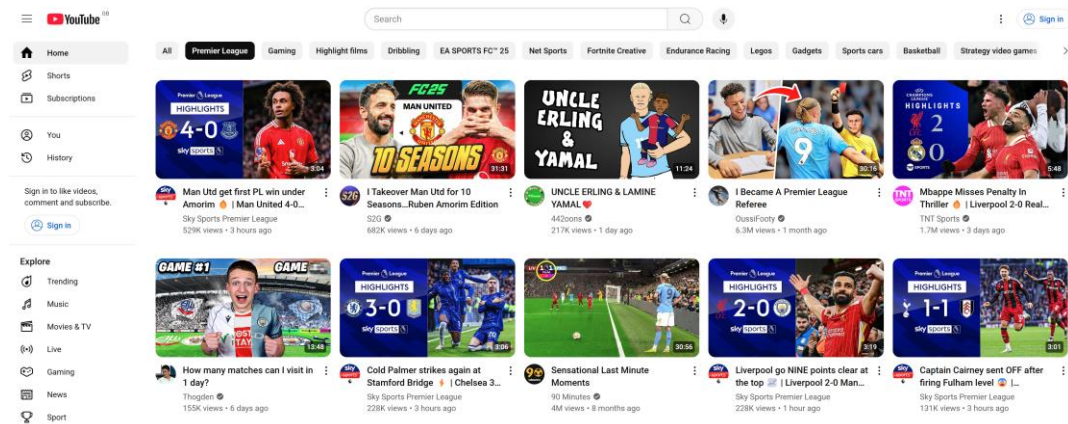
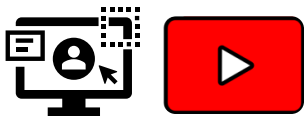


Linear Contextual Bandits


- T rounds, K arms ()

For $t = 1, \dots, T$:

1) **Algorithm** observes arm-specific contexts $x_{t,1}^* = \quad, \dots, x_{t,K}^* = \quad \in \mathbb{R}^d$



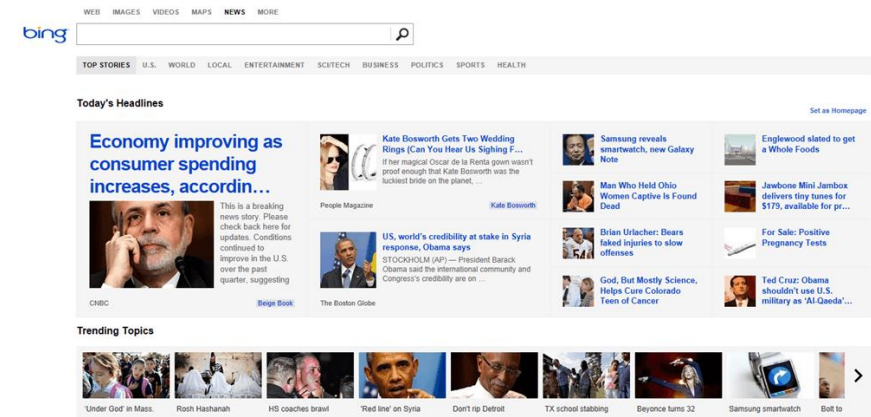
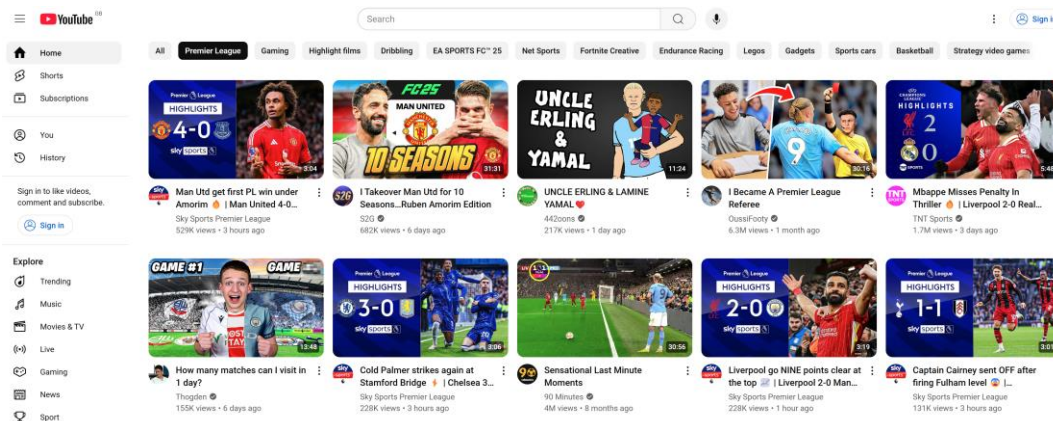
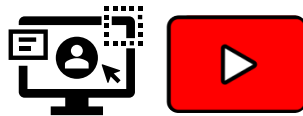
Linear Contextual Bandits

- T rounds, K arms ()


For $t = 1, \dots, T$:

1) **Algorithm** observes **arm-specific** contexts $x_{t,1}^* =$, ... , $x_{t,K}^* =$ $\in \mathbb{R}^d$

2) **Algorithm** plays arm $i_t =$ $\in [K]$ and receives reward $r_t^*(i_t) :=$

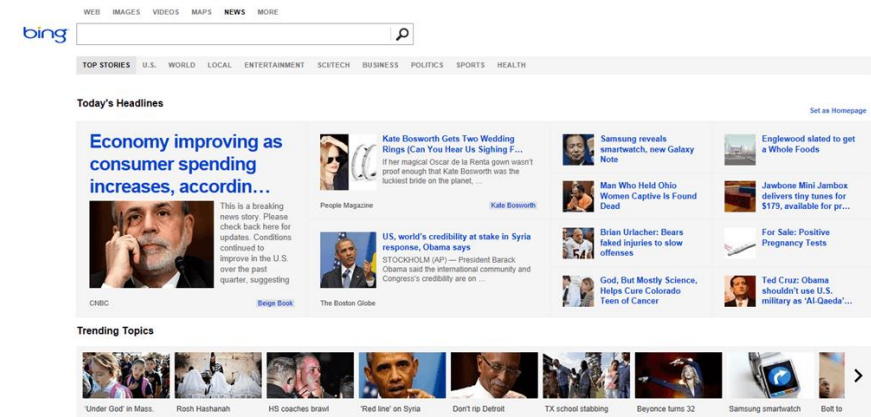
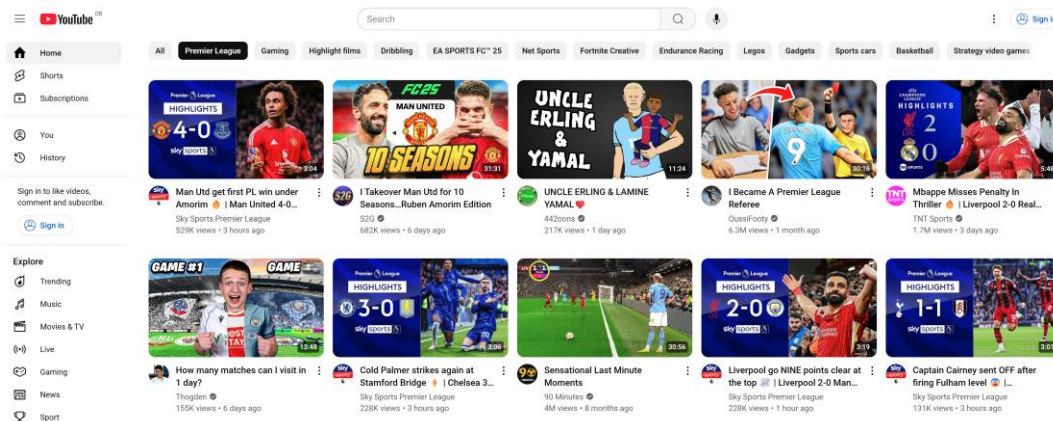
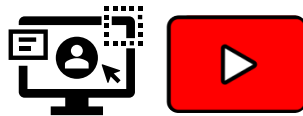


Linear Contextual Bandits


- T rounds, K arms ()

For $t = 1, \dots, T$:

- 1) **Algorithm** observes **arm-specific** contexts $x_{t,1} = \text{YouTube icon}$, \dots , $x_{t,K} = \text{YouTube icon} \in \mathbb{R}^d$
- 2) **Algorithm** plays arm $i_t = \dots \in [K]$ and receives reward $r_t^*(i_t) := \dots$



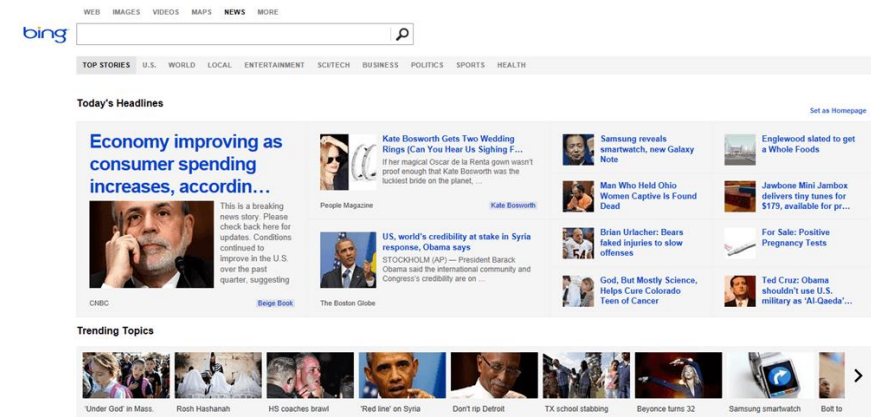
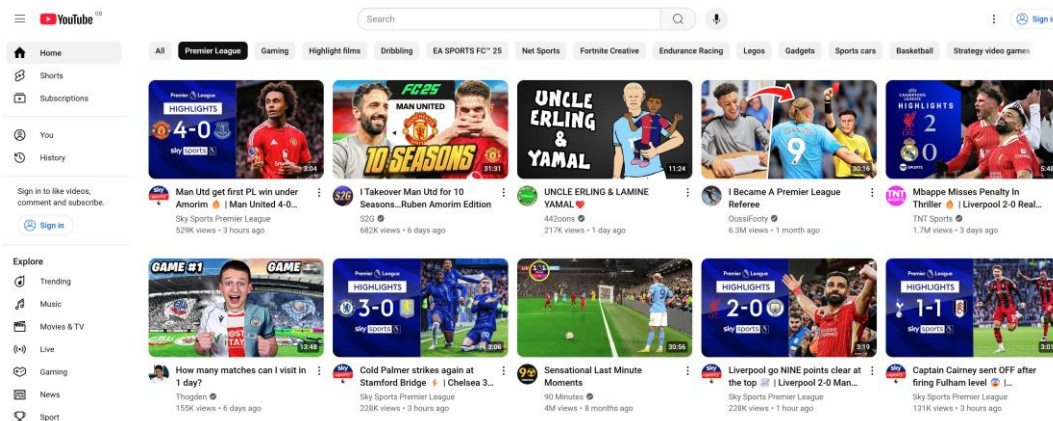
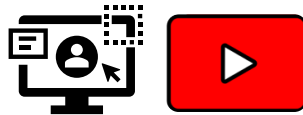
Linear Contextual Bandits

- T rounds, K arms ()


For $t = 1, \dots, T$:

1) **Algorithm** observes **arm-specific** contexts $x_{t,1}^* = \text{red square with play button and painter icon}$, \dots , $x_{t,K}^* = \text{red square with play button and painter icon} \in \mathbb{R}^d$

2) **Algorithm** plays arm $i_t = \text{red square with painter icon} \in [K]$ and receives reward $r_t^*(i_t) :=$



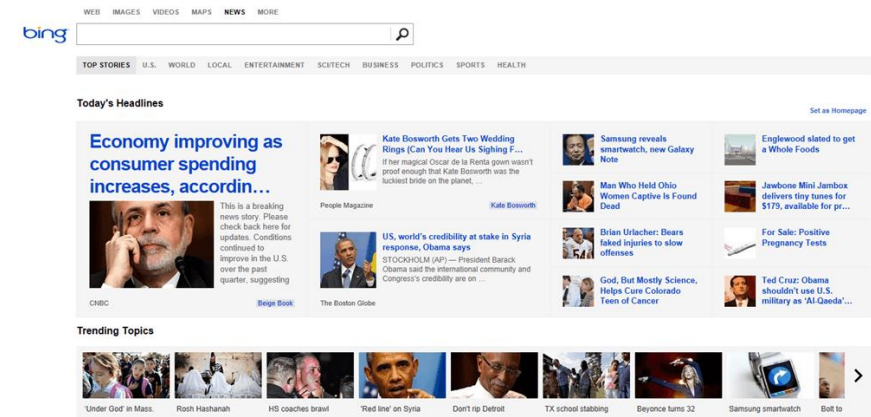
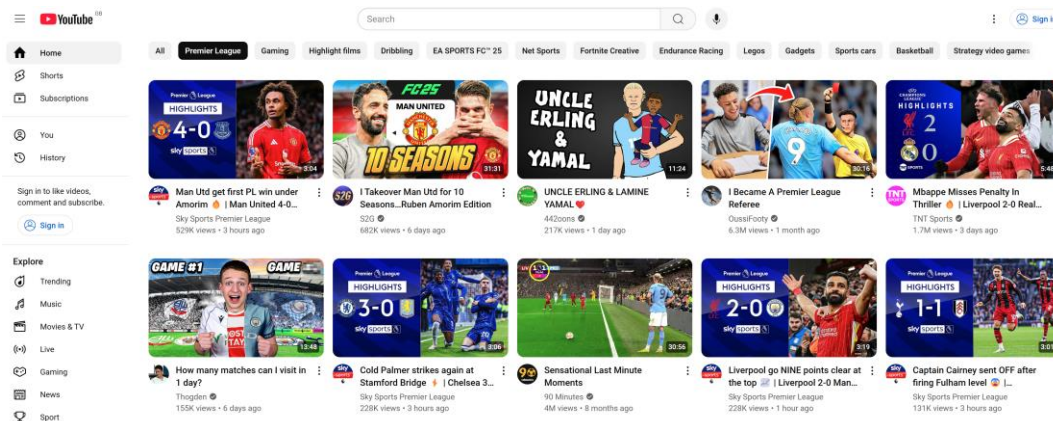
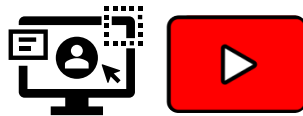
Linear Contextual Bandits

- T rounds, K arms ()

For $t = 1, \dots, T$:

1) **Algorithm** observes arm-specific contexts $x_{t,1}^* = \text{play button icon}$, \dots , $x_{t,K}^* = \text{painter icon}$ $\in \mathbb{R}^d$

2) **Algorithm** plays arm $i_t = \text{painter icon}$ $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$ 




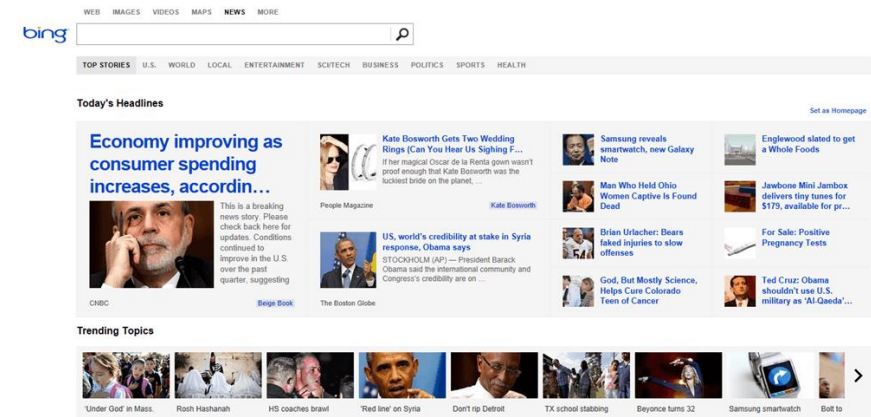
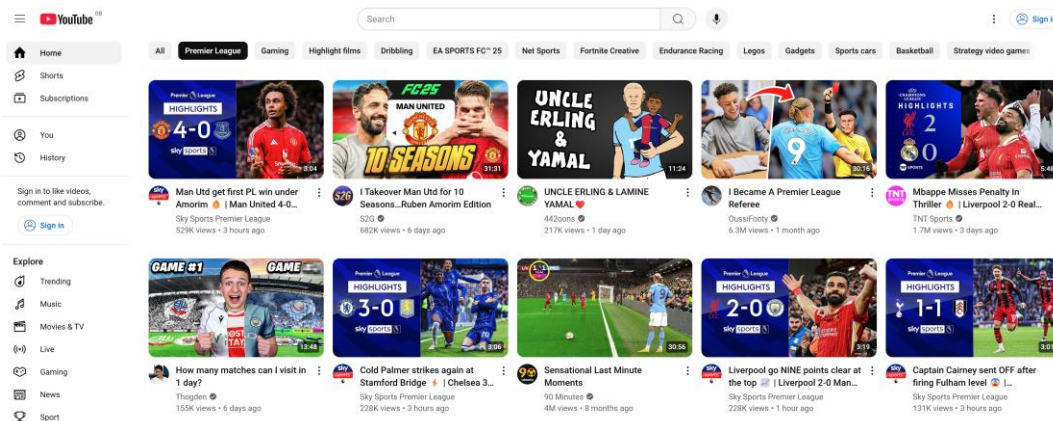
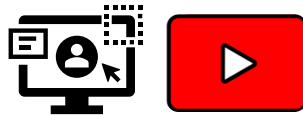
Linear Contextual Bandits

- T rounds, K arms ()


For $t = 1, \dots, T$:

1) **Algorithm** observes arm-specific contexts $x_{t,1}^* = \text{YouTube icon}$, \dots , $x_{t,K}^* = \text{YouTube icon} \in \mathbb{R}^d$


2) **Algorithm** plays arm $i_t = \text{YouTube icon}$ $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$ 
↑
unknown




Linear Contextual Bandits


- T rounds, K arms ()

For $t = 1, \dots, T$:


- 1) **Algorithm** observes **arm-specific** contexts $x_{t,1}^* = \text{red square with play button and painter icon}$, \dots , $x_{t,K}^* = \text{red square with play button and painter icon}$ $\in \mathbb{R}^d$
- 2) **Algorithm** plays arm $i_t = \text{red square with painter icon}$ $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$ 



unknown

Linear Contextual Bandits

- T rounds, K arms ()

For $t = 1, \dots, T$:


- 1) **Algorithm** observes **arm-specific** contexts $x_{t,1}^* = \text{red square with play button}$, \dots , $x_{t,K}^* = \text{red square with play button}$ $\in \mathbb{R}^d$
- 2) **Algorithm** plays arm $i_t = \text{red square with painter icon}$ $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$ 


unknown

Algorithm maximizes cumulative reward


$$\sum_{t=1}^T r_t(i_t)$$

Linear Contextual Bandits

- T rounds, K arms ()

For $t = 1, \dots, T$:


1) **Algorithm** observes **arm-specific** contexts $x_{t,1}^* = \text{play button icon}$, \dots , $x_{t,K}^* = \text{artist icon}$ $\in \mathbb{R}^d$

2) **Algorithm** plays arm $i_t = \text{artist icon}$ $\in [K]$ and receives reward $r_t^*(i_t) := \underbrace{\langle \theta^*, x_{t,i_t}^* \rangle}_{\text{unknown}} + \eta_t$ 

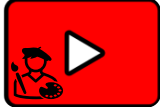




Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Linear Contextual Bandits

- T rounds, K arms ()

For $t = 1, \dots, T$:

- 1) **Algorithm** observes **arm-specific** contexts $x_{t,1}^* =$ , \dots , $x_{t,K}^* =$  $\in \mathbb{R}^d$
private information 
- 2) **Algorithm** plays arm $i_t =$  $\in [K]$ and receives reward $r_t^*(i_t) := \underbrace{\langle \theta^*, x_{t,i_t}^* \rangle}_{\text{unknown}} + \eta_t$ 

Algorithm minimizes expected regret



$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$


Linear Contextual Bandits


- T rounds, K arms ()

For $t = 1, \dots, T$:

1) **Algorithm** ~~observes~~ **arm-specific** contexts $x_{t,1}^* =$ , \dots , $x_{t,K}^* =$  $\in \mathbb{R}^d$

2) **Algorithm** plays arm $i_t =$  $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$ 

private information 

unknown 

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Strategic Linear Contextual Bandits

arm = strategic agent



For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \quad \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \quad \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \quad \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Strategic Linear Contextual Bandits

arm = strategic agent



For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{[arm icon]} \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{[arm icon]} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{[arm icon]} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Strategic Linear Contextual Bandits

arm = strategic agent



For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{[red box with artist icon and play button]} \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{[red box with artist icon and rainbow triangle]} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{[red box with artist icon]} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Strategic Linear Contextual Bandits

arm = strategic agent



For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{[arm icon with play button]} \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{[arm icon with rainbow triangle]} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{[arm icon]} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Every Arm i maximizes its **#selections**

$$\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$$

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤🎨} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{👤🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$




Every **Arm** i maximizes its **#selections**

$$\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$$

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$: \leftarrow *repeated interaction* 

- 1) Every arm $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{} \in \mathbb{R}^d$
- 2) Every arm $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$


Every Arm i maximizes its **#selections**

$$\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$$

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$: \leftarrow *repeated interaction* 

- 1) Every arm $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
- 2) Every arm $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤🎨} \in \mathbb{R}^d$ to the **Algorithm** *unbounded manipulation* 
- 3) **Algorithm** plays arm $i_t = \text{👨🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$



Every Arm i maximizes its **#selections**

$$\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$$

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$: *repeated interaction* 

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \img alt="red box with play button icon" data-bbox="652 253 713 331" $\in \mathbb{R}^d$$
 - 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \img alt="red box with play button and paint palette icon" data-bbox="605 351 667 429" $\in \mathbb{R}^d$ to the **Algorithm** *unbounded manipulation* $
 - 3) **Algorithm** plays arm $i_t = \img alt="red box with artist icon" data-bbox="357 444 421 524" $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$$
- multiple competing agents* 

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$



Every **Arm** i maximizes its **#selections**

$$\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$$

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$


For $t = 1, \dots, T$: *repeated interaction* 

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \img alt="red box with play button icon" data-bbox="652 254 713 331" $\in \mathbb{R}^d$$
 - 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \img alt="red box with play button and paint palette icon" data-bbox="604 351 667 428" $\in \mathbb{R}^d$ to the **Algorithm** *unbounded manipulation* $
 - 3) **Algorithm** plays arm $i_t = \img alt="red box with artist icon" data-bbox="357 444 421 524" $\in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$$
- multiple competing agents* 

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Every **Arm** i maximizes its **#selections**




 $\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$

far-sighted agents

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$


For $t = 1, \dots, T$: *repeated interaction* 

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
 - 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤🎮} \in \mathbb{R}^d$ to the **Algorithm** *unbounded manipulation* 
 - 3) **Algorithm** plays arm $i_t = \text{👨🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$
- multiple competing agents*  *unknown environment* 

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Every **Arm** i maximizes its **#selections**




 $\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$

far-sighted agents

Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$


For $t = 1, \dots, T$: *repeated interaction* 

- 1) Every arm $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
- 2) Every arm $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤🎮} \in \mathbb{R}^d$ to the **Algorithm** *unbounded manipulation* 
- 3) **Algorithm** plays arm $i_t = \text{👨🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$
multiple competing agents  *unknown environment* 

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Every Arm i maximizes its **#selections**




 $\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$
far-sighted agents



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$


For $t = 1, \dots, T$: *repeated interaction* 

- 1) Every arm $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
- 2) Every arm $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤🎮} \in \mathbb{R}^d$ to the **Algorithm** *unbounded manipulation* 
- 3) **Algorithm** plays arm $i_t = \text{👨🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$
multiple competing agents  *known environment* 

Algorithm minimizes expected regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \max_{i \in [K]} \langle \theta^*, x_{t,i}^* \rangle - \langle \theta^*, x_{t,i_t}^* \rangle \right]$$

Every Arm i maximizes its **#selections**

 $\mathbb{E} \left[\sum_{t=1}^T 1(i_t = i) \right]$
far-sighted agents



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤🎨} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{👨🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes its context $x_{t,i}^* = \text{👤▶} \in \mathbb{R}^d$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{👤📺} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{👨🎨} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

\nearrow
known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$



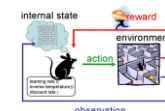
Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes **true avg reward** $\mu_{t,i}^* = \text{rock icon} \in \mathbb{R}$
- 2) Every **arm** $i \in [K]$ reports a **gamed** context $x_{t,i} = \text{TV icon} \in \mathbb{R}^d$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{artist icon} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

\nearrow
known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

- 1) Every **arm** $i \in [K]$ **privately** observes **true avg reward** $\mu_{t,i}^* = \text{rock} \in \mathbb{R}$
- 2) Every **arm** $i \in [K]$ reports a **gamed value** $\mu_{t,i} = \text{diamond} \in \mathbb{R}$ to the **Algorithm**
- 3) **Algorithm** plays arm $i_t = \text{artist} \in [K]$ and receives reward $r_t^*(i_t) := \langle \theta^*, x_{t,i_t}^* \rangle + \eta_t$

\nearrow
known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

1) Every **arm** $i \in [K]$ **privately** observes **true avg reward** $\mu_{t,i}^* = \text{rock} \in \mathbb{R}$

2) Every **arm** $i \in [K]$ reports a **gamed value** $\mu_{t,i} = \text{diamond} \in \mathbb{R}$ to the **Algorithm**

3) **Algorithm** plays arm $i_t = \text{artist} \in [K]$ and receives reward $r_t^*(i_t) := \mu_{t,i_t}^* + \eta_t$

known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

1) Every **arm** $i \in [K]$ **privately** observes **true avg reward** $\mu_{t,i}^* = \text{rock} \in \mathbb{R}$

2) Every **arm** $i \in [K]$ reports a **gamed value** $\mu_{t,i} = \text{diamond} \in \mathbb{R}$ to the **Algorithm**

3) **Algorithm** plays arm $i_t = \text{artist} \in [K]$ and receives reward $r_t^*(i_t) := \mu_{t,i_t}^* + \eta_t$

known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$
- We're happy when the arms "cooperate" and the **reports** $\mu_{t,i}$ match the **truth** $\mu_{t,i}^*$



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

1) Every **arm** $i \in [K]$ **privately** observes **true avg reward** $\mu_{t,i}^* = \text{rock} \in \mathbb{R}$

2) Every **arm** $i \in [K]$ reports a **gamed value** $\mu_{t,i} = \text{diamond} \in \mathbb{R}$ to the **Algorithm**

3) **Algorithm** plays arm $i_t = \text{artist} \in [K]$ and receives reward $r_t^*(i_t) := \mu_{t,i_t}^* + \eta_t$

known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$

- We're happy when the arms "cooperate" and the **reports** $\mu_{t,i}$ **diamond** match the **truth** $\mu_{t,i}^*$ **rock**



Strategic Linear Contextual Bandits

Arms respond in **Equilibrium**:
arm strategies $\in \text{NE}(\text{Algorithm})$

For $t = 1, \dots, T$:

1) Every **arm** $i \in [K]$ **privately** observes **true avg reward** $\mu_{t,i}^* = \text{rock} \in \mathbb{R}$

2) Every **arm** $i \in [K]$ reports a **gamed value** $\mu_{t,i} = \text{diamond} \in \mathbb{R}$ to the **Algorithm**

3) **Algorithm** plays arm $i_t = \text{artist} \in [K]$ and receives reward $r_t^*(i_t) := \mu_{t,i_t}^* + \eta_t$

known environment



- Let's swap features $x \in \mathbb{R}^d$ for the implied **avg reward** $\mu := \langle \theta^*, x \rangle \in \mathbb{R}$
- We're happy when the arms "cooperate" and the **reports** $\mu_{t,i}$ **diamond** match the **truth** $\mu_{t,i}^*$ **rock**

*How can we incentivize the arms to "cooperate" with us
while minimizing regret?*



Grim Trigger Strategy

How can we incentivize the arms to “cooperate” with us?

Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$
- 2) each arm tells us $\mu_{t,i}$
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$

Grim Trigger Strategy

How can we incentivize the arms to “cooperate” with us?

Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$
- 2) each arm tells us $\mu_{t,i}$
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$

Grim Trigger Strategy



How can we incentivize the arms to “cooperate” with us?

From Iterated Social Dilemmas: *“If you defect once, I will defect permanently.”*

Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$
- 2) each arm tells us $\mu_{t,i}$
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$

Grim Trigger Strategy



How can we incentivize the arms to “cooperate” with us?

From Iterated Social Dilemmas: *“If you defect once, I will defect permanently.”*

*If **reported avg reward** $\mu_{t,i}$ $>$ **true avg reward** $\mu_{t,i}^*$,*

eliminate arm i forever.



Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$*
- 2) each arm tells us $\mu_{t,i}$*
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$*

Grim Trigger Strategy



How can we incentivize the arms to “cooperate” with us?

From Iterated Social Dilemmas: *“If you defect once, I will defect permanently.”*

*If **reported avg reward** $\mu_{t,i}$ $>$ **true avg reward** $\mu_{t,i}^*$,*

eliminate arm i forever.



But we don't observe $\mu_{t,i}^$*



Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$*
- 2) each arm tells us $\mu_{t,i}$*
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$*

Grim Trigger Strategy



How can we incentivize the arms to “cooperate” with us?

From Iterated Social Dilemmas: *“If you defect once, I will defect permanently.”*

*If **reported avg reward** $\mu_{t,i}$ > **observed (true) reward** $r_{t,i}^*$,*

eliminate arm i forever.



But we don't observe $\mu_{t,i}^$*



Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$*
- 2) each arm tells us $\mu_{t,i}$*
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$*

Grim Trigger Strategy



How can we incentivize the arms to “cooperate” with us?

From Iterated Social Dilemmas: *“If you defect once, I will defect permanently.”*

*If **reported avg reward** $\mu_{t,i}$ **>** **observed (true) reward** $r_{t,i}^*$,*

eliminate arm i forever.



But $r_{t,i}^$ is noisy*



Repeatedly:

- 1) each arm privately observes $\mu_{t,i}^*$*
- 2) each arm tells us $\mu_{t,i}$*
- 3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$*

Grim Trigger Strategy

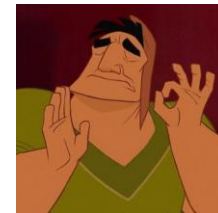


How can we incentivize the arms to “cooperate” with us?

From Iterated Social Dilemmas: *“If you defect once, I will defect permanently.”*

*If **reported avg reward** $\sum_t \mu_{t,i}$ > **optimistic estimate** $\sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$*

eliminate arm i forever.



Repeatedly:

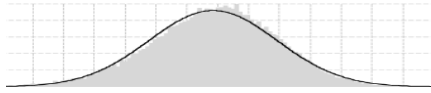
1) each arm privately observes $\mu_{t,i}^*$

2) each arm tells us $\mu_{t,i}$

3) We observe $r_{t,i_t}^* \sim D(\mu_{t,i_t}^*)$

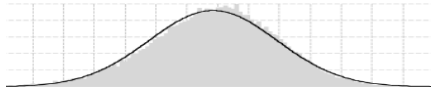
What does the Grim Trigger mean for the arms?

What does the Grim Trigger mean for the arms?



The observations $\sum_t r_{t,i}^*$ *concentrate* around the average $\sum_t \mu_{t,i}^*$

What does the Grim Trigger mean for the arms?



The observations $\sum_t r_{t,i}^*$ *concentrate* around the average $\sum_t \mu_{t,i}^*$

This implies that with high probability:

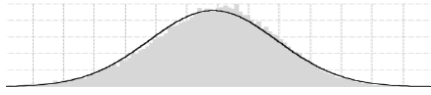
If $\sum_t \mu_{t,i} - \sum_t \mu_{t,i}^* > 2\sqrt{n_\tau(i)}$, we eliminate arm i .

reports — *truth*



What does the Grim Trigger mean for the arms?

The observations $\sum_t r_{t,i}^*$ *concentrate* around the average $\sum_t \mu_{t,i}^*$



This implies that with high probability:

If $\sum_t \mu_{t,i} - \sum_t \mu_{t,i}^* > 2\sqrt{n_\tau(i)}$, we eliminate arm i .

reports — *truth*



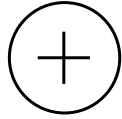
The arms can misreport to us ... but not too often !

Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Selection:

Play $i_t = \arg \max_{i \in \text{alive}} \mu_{t,i}$



Grim Trigger:

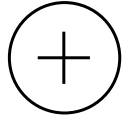


If $\sum_t \mu_{t,i} > \sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$, eliminate arm i .

Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Selection:

Play $i_t = \arg \max_{i \in \text{alive}} \mu_{t,i}$



Grim Trigger:



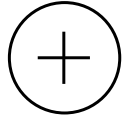
If $\sum_t \mu_{t,i} > \sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$, eliminate arm i .

Theoretical Results:

Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Selection:

Play $i_t = \arg \max_{i \in \text{alive}} \mu_{t,i}$



Grim Trigger:



If $\sum_t \mu_{t,i} > \sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$, eliminate arm i .

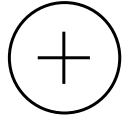
Theoretical Results:

- 1) Under **GGTM**, being **truthful** is a \sqrt{T} -Nash Equilibrium.

Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Selection:

Play $i_t = \arg \max_{i \in \text{alive}} \mu_{t,i}$



Grim Trigger:



If $\sum_t \mu_{t,i} > \sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$, eliminate arm i .

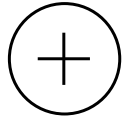
Theoretical Results:

- 1) Under **GGTM**, being **truthful** is a \sqrt{T} -Nash Equilibrium.
- 2) If the arms play any NE under **GGTM**, then $R_T(\mathbf{GGTM}) \leq K\sqrt{T} + K^2\sqrt{KT}$

Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Selection:

Play $i_t = \arg \max_{i \in \text{alive}} \mu_{t,i}$



Grim Trigger:

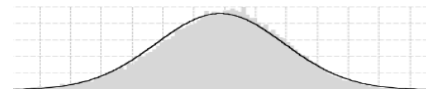


If $\sum_t \mu_{t,i} > \sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$, eliminate arm i .

Theoretical Results:

- 1) Under **GGTM**, being **truthful** is a \sqrt{T} -Nash Equilibrium.
- 2) If the arms play any NE under **GGTM**, then $R_T(\text{GGTM}) \leq K\sqrt{T} + K^2\sqrt{KT}$

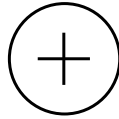
*cost of manipulation
& uncertainty*



Greedy Grim Trigger Mechanism (**GGTM**)

Greedy Selection:

Play $i_t = \arg \max_{i \in \text{alive}} \mu_{t,i}$



Grim Trigger:



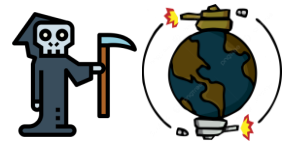
If $\sum_t \mu_{t,i} > \sum_t r_{t,i}^* + \sqrt{n_\tau(i)}$, eliminate arm i .

Theoretical Results:

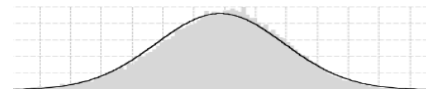
1) Under **GGTM**, being **truthful** is a \sqrt{T} -Nash Equilibrium.

2) If the arms play any NE under **GGTM**, then $R_T(\text{GGTM}) \leq K\sqrt{T} + K^2\sqrt{KT}$

cost of mechanism design



*cost of manipulation
& uncertainty*



Suppose the environment θ^* is **unknown** ...

Things get complicated ...

We observe: gamed context $x_{t,i}$ and reward $r_{t,i}^* := \langle \theta^*, x_{t,i}^* \rangle + \eta_t$

We don't observe: true context $x_{t,i}^*$ and parameter θ^*

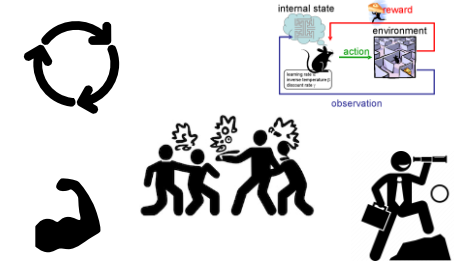
The arms can manipulate our estimate of θ^* ...





Estimating θ^* accurately becomes **impossible**?!

Another time ...

Short Recap

- Strategic Interactive Decision-Making
 - **Reinforcement Learning** + **Mechanism Design**
 - Objective: **Strategic Robustness** + **Incentive Alignment**



- Strategic Linear Contextual Bandits
 - Strategic agents  manipulating contexts 
 - Grim Trigger Mechanism  from Iterated Social Dilemmas
 - Mechanism Design becomes *approximate* 

- *There are many more problems like this left to study ...*

