

# Nonstochastic Bandit Problems on Graphs

Nicolò Cesa-Bianchi

Università degli Studi di Milano



# Theory of repeated games



James Hannan  
(1922–2010)



David Blackwell  
(1919–2010)

Learning to play a game (1956)

Play a game repeatedly against a possibly suboptimal opponent

# Zero-sum 2-person games played more than once

	1	2	...	M
1	$\ell(1,1)$	$\ell(1,2)$	...	
2	$\ell(2,1)$	$\ell(2,2)$	...	
$\vdots$	$\vdots$	$\vdots$	$\ddots$	
K				

$K \times M$  known loss matrix over  $\mathbb{R}$

- Row player (**player**) has  $K$  actions
- Column player (**opponent**) has  $M$  actions

For each game round  $t = 1, 2, \dots$

- Player chooses action  $i_t$  and opponent chooses action  $y_t$
- The player suffers loss  $\ell(i_t, y_t)$  (= gain of opponent)

Player can learn from opponent's history of past choices  $y_1, \dots, y_{t-1}$



# Prediction with expert advice



Volodya Vovk



Manfred Warmuth

	$t = 1$	$t = 2$	$\dots$
1	$\ell_1(1)$	$\ell_2(1)$	$\dots$
2	$\ell_1(2)$	$\ell_2(2)$	$\dots$
$\vdots$	$\vdots$	$\vdots$	$\ddots$
K	$\ell_1(K)$	$\ell_2(K)$	

Play an unknown loss matrix

Opponent's moves  $y_1, y_2, \dots$  define a **sequential prediction problem** with a **time-varying loss** function  $\ell(i_t, y_t) = \ell_t(i_t)$



# Playing the experts game

## A sequential decision problem

- $K$  actions
- Unknown deterministic assignment of losses to actions  
 $\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K$  for  $t = 1, 2, \dots$



For  $t = 1, 2, \dots$

# Playing the experts game

## A sequential decision problem

- $K$  actions
- Unknown deterministic assignment of losses to actions  
 $\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K$  for  $t = 1, 2, \dots$



For  $t = 1, 2, \dots$

- 1 Player picks an action  $I_t$  (possibly using randomization) and incurs loss  $\ell_t(I_t)$

# Playing the experts game

## A sequential decision problem

- $K$  actions
- Unknown deterministic assignment of losses to actions  
 $\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K$  for  $t = 1, 2, \dots$



For  $t = 1, 2, \dots$

- 1 Player picks an action  $I_t$  (possibly using randomization) and incurs loss  $\ell_t(I_t)$
- 2 Player gets **feedback information**:  $\ell_t(1), \dots, \ell_t(K)$

# Regret analysis

## Regret

$$R_T \stackrel{\text{def}}{=} \mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t) \right] - \min_{i=1, \dots, K} \sum_{t=1}^T \ell_t(i) \stackrel{\text{want}}{=} o(T)$$





# Regret analysis

## Regret

$$R_T \stackrel{\text{def}}{=} \mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t) \right] - \min_{i=1, \dots, K} \sum_{t=1}^T \ell_t(i) \stackrel{\text{want}}{=} o(T)$$

## Lower bound using random losses

[Experts' paper, 1997]

- $\ell_t(i) \rightarrow L_t(i) \in \{0, 1\}$  independent random coin flip

- For any player strategy  $\mathbb{E} \left[ \sum_{t=1}^T L_t(I_t) \right] = \frac{T}{2}$

- Then the expected regret is

$$\mathbb{E} \left[ \max_{i=1, \dots, K} \sum_{t=1}^T \left( \frac{1}{2} - L_t(i) \right) \right] = (1 - o(1)) \sqrt{\frac{T \ln K}{2}}$$

for  $K, T \rightarrow \infty$

# Exponentially weighted forecaster (Hedge)

At time  $t$  pick action  $I_t = i$  with probability proportional to

$$\exp\left(-\eta \sum_{s=1}^{t-1} \ell_s(i)\right)$$

the sum at the exponent is the **total loss** of action  $i$  up to now

Regret bound

[Experts' paper, 1997]

- If  $\eta = \sqrt{\frac{\ln K}{8T}}$  then  $R_T \leq \sqrt{\frac{T \ln K}{2}}$
- This matches the asymptotic lower bound, **including constants**



# The bandit problem: playing an unknown game

- $K$  actions
- Unknown deterministic assignment of losses to actions  
 $\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K$  for  $t = 1, 2, \dots$



For  $t = 1, 2, \dots$



# The bandit problem: playing an unknown game

- $K$  actions
- Unknown deterministic assignment of losses to actions  
 $\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K$  for  $t = 1, 2, \dots$



For  $t = 1, 2, \dots$

- 1 Player picks an action  $I_t$  (possibly using randomization) and incurs loss  $\ell_t(I_t)$



# The bandit problem: playing an unknown game

- $K$  actions
- Unknown deterministic assignment of losses to actions  
 $\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K$  for  $t = 1, 2, \dots$



For  $t = 1, 2, \dots$

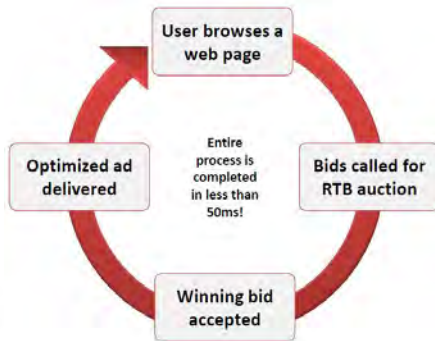
- 1 Player picks an action  $I_t$  (possibly using randomization) and incurs loss  $\ell_t(I_t)$
- 2 Player gets **feedback information**: Only  $\ell_t(I_t)$  is revealed



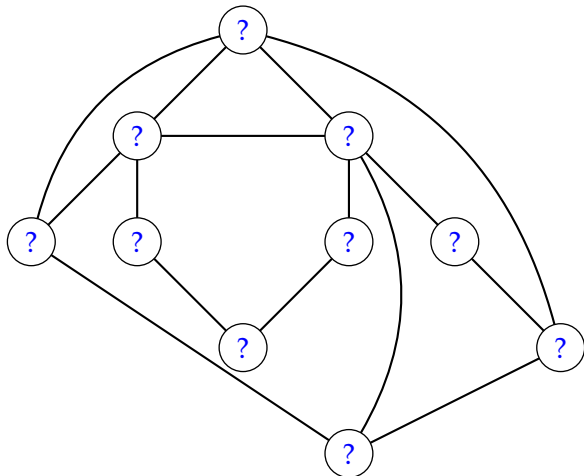
# A growing range of applications

- Ad placement
- Dynamic content/layout optimization
- Recommender systems
- Clinical trials
- Network protocol optimization



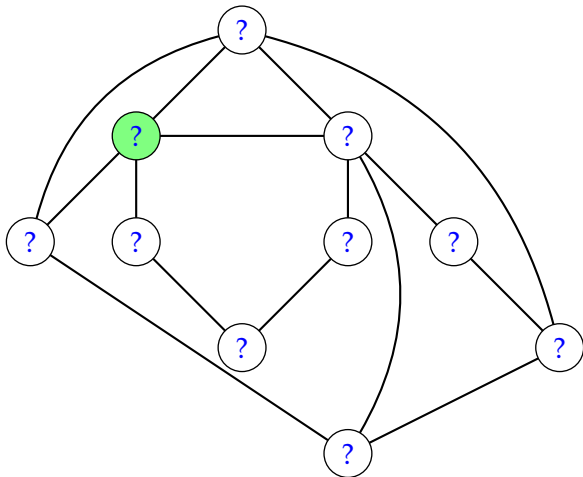


# A graph of relationships over actions

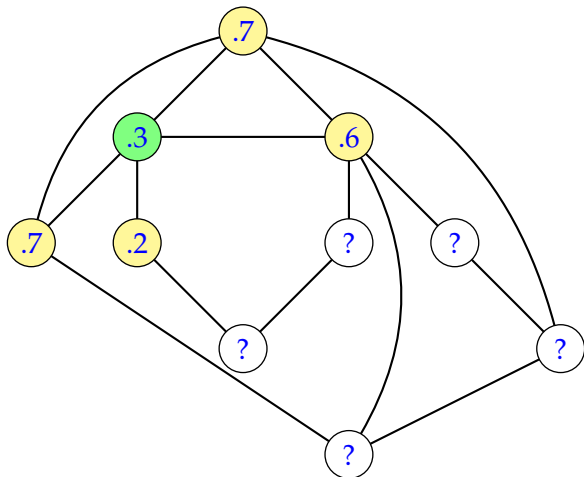




# A graph of relationships over actions

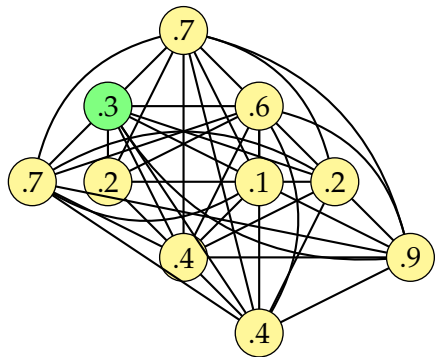


# A graph of relationships over actions

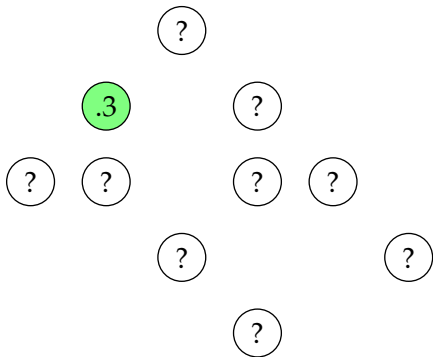


# Recovering expert and bandit settings

Experts: clique



Bandits: empty graph



# Hedge revisited

Player's strategy [Alon, C-B, Gentile, Mannor, Mansour and Shamir, 2013]

- $\mathbb{P}_t(I_t = i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_s(i)\right) \quad i = 1, \dots, K$
- $\hat{\ell}_t(i) = \begin{cases} \frac{\ell_t(i)}{\mathbb{P}_t(\ell_t(i) \text{ observed})} & \text{if } \ell_t(i) \text{ is observed} \\ 0 & \text{otherwise} \end{cases}$



# Hedge revisited

Player's strategy [Alon, C-B, Gentile, Mannor, Mansour and Shamir, 2013]

- $\mathbb{P}_t(I_t = i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_s(i)\right) \quad i = 1, \dots, K$
- $\hat{\ell}_t(i) = \begin{cases} \frac{\ell_t(i)}{\mathbb{P}_t(\ell_t(i) \text{ observed})} & \text{if } \ell_t(i) \text{ is observed} \\ 0 & \text{otherwise} \end{cases}$

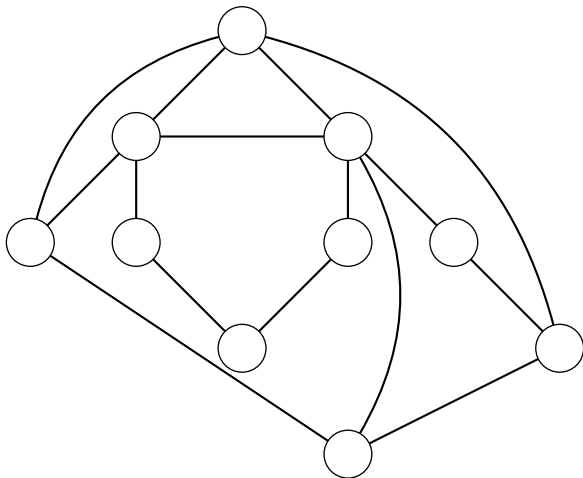
Importance sampling estimator

$$\mathbb{E}_t[\hat{\ell}_t(i)] = \ell_t(i) \quad \text{unbiasedness}$$
$$\mathbb{E}_t[\hat{\ell}_t(i)^2] \leq \frac{1}{\mathbb{P}_t(\ell_t(i) \text{ observed})} \quad \text{variance control}$$



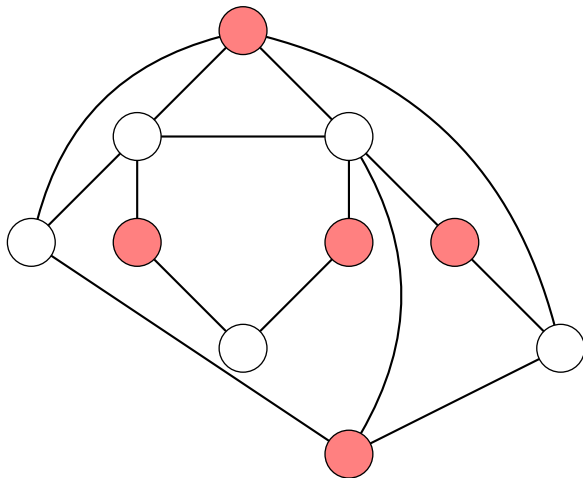
# Independence number $\alpha(G)$

The size of the largest **independent set**



# Independence number $\alpha(G)$

The size of the largest **independent set**



# Regret bound

$$\begin{aligned} R_T &\leq \frac{\ln K}{\eta} + \frac{\eta}{2} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^K \mathbb{P}_t(I_t = i) \mathbb{E}_t \left[ \widehat{\ell}_t(i)^2 \right] \right] \\ &\leq \frac{\ln K}{\eta} + \frac{\eta}{2} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^K \frac{\mathbb{P}_t(I_t = i)}{\mathbb{P}_t(\ell_t(i) \text{ is observed})} \right] \quad \text{variance control} \\ &\leq \frac{\ln K}{\eta} + \frac{\eta}{2} T \alpha(G) = \sqrt{T \alpha(G) \ln K} \end{aligned}$$

$\alpha(G)$  is the **independence number** of  $G$





# Regret bound

$$\begin{aligned} R_T &\leq \frac{\ln K}{\eta} + \frac{\eta}{2} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^K \mathbb{P}_t(I_t = i) \mathbb{E}_t \left[ \widehat{\ell}_t(i)^2 \right] \right] \\ &\leq \frac{\ln K}{\eta} + \frac{\eta}{2} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^K \frac{\mathbb{P}_t(I_t = i)}{\mathbb{P}_t(\ell_t(i) \text{ is observed})} \right] \quad \text{variance control} \\ &\leq \frac{\ln K}{\eta} + \frac{\eta}{2} T \alpha(G) = \sqrt{T \alpha(G) \ln K} \end{aligned}$$

$\alpha(G)$  is the **independence number** of  $G$

## Special cases

**Experts** (clique):

$$\alpha(G) = 1$$

$$R_T \leq \sqrt{T \ln K}$$

**Bandits** (empty graph):

$$\alpha(G) = K$$

$$R_T \leq \sqrt{TK \ln K}$$

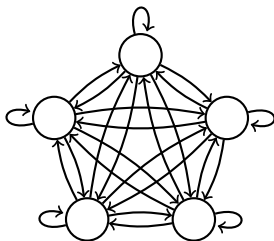
## Directed



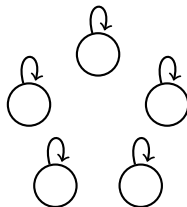
## Interventions



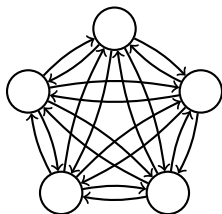
# Old and new examples



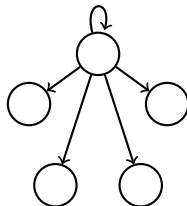
Experts



Bandits



Cops & Robbers



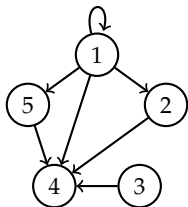
Revealing Action



# A characterization of feedback graphs

A vertex of  $G$  is:

- **observable** if it has at least one incoming edge (possibly a self-loop)
- **strongly observable** if it has either a self-loop or incoming edges from all other vertices
- **weakly observable** if it is observable but not strongly observable



- 3 is not observable
- 2 and 5 are weakly observable
- 1 and 4 are strongly observable



# Minimax rates

Hedge uses additional exploration due to reduced feedback

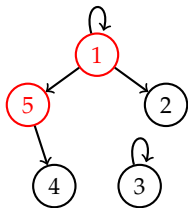
$G$  is **strongly observable**  $R_T = \tilde{\Theta}\left(\sqrt{\alpha(G)T}\right)$

Hedge mixed with uniform distribution

$G$  is **weakly observable**  $R_T = \tilde{\Theta}\left(T^{2/3}\delta(G)\right)$

Hedge mixed with uniform distribution on a weakly dominating set

$G$  is **not observable**  $R_T = \Theta(T)$

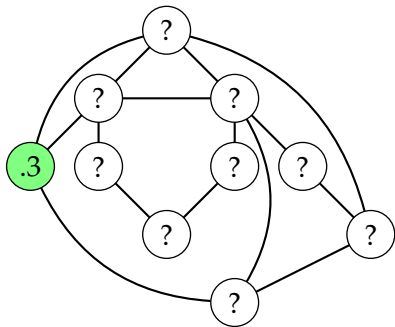


Weakly dominating set

$\delta(G)$  is the size of the smallest set that dominates all weakly observable nodes of  $G$

- **Back in the bandit model:** we only observe the loss of the chosen action
- But we want to take advantage of scenarios where  $\ell_t$  is **smooth** over the graph:

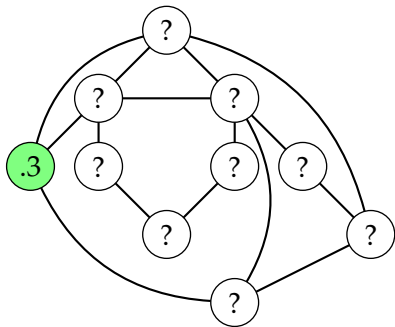
$$\sum_{(i,j) \in E} (\ell_t(i) - \ell_t(j))^2 \text{ is small}$$



- **Back in the bandit model:** we only observe the loss of the chosen action
- But we want to take advantage of scenarios where  $\ell_t$  is **smooth** over the graph:

$$\sum_{(i,j) \in E} (\ell_t(i) - \ell_t(j))^2 \text{ is small}$$

- We also need to know the **smallest loss**  $\min_{i=1,\dots,K} \ell_t(i)$   
(no need to know the action achieving it)



# Regret with smooth losses

$$R_T \leq \sqrt{(\ln K) \frac{C^2}{\lambda_2(L)} T}$$

for the Laplacian  $L$  of any simple and connected graph such that

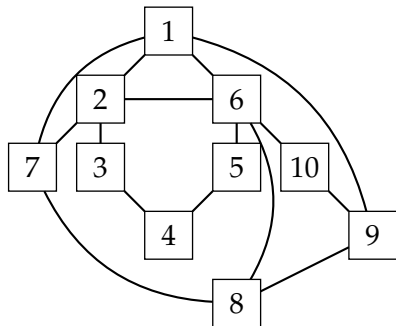
$$\underbrace{\sum_{(i,j) \in E} (\ell_t(i) - \ell_t(j))^2}_{\ell_t^T L \ell_t} \leq C^2 \quad \text{for all } t = 1, \dots, T$$

- $\lambda_2(K) \in [0, K]$  is the smallest nonzero eigenvalue of  $L$  measuring connectivity of the graph
- This regret bound vanishes when  $\ell_t(1) = \dots = \ell_t(K)$  for all  $t$





- $N$  agents sitting on the vertices of an unknown **communication graph**
- Agents **cooperate** to solve a common bandit problem
- Each agent runs an instance of the same bandit algorithm



# The communication protocol with delay $d$

For each  $t = 1, \dots, T$  each agent  $v \in V$  does the following:

- 1 Plays an action  $I_t(v)$  drawn according to his private distribution  $p_t(v)$  observing loss  $\ell_t(I_t(v))$  (same loss vector for all agents)
- 2 Sends to his neighbors the message

$$m_t(v) = \langle t, v, I_t(v), \ell_t(I_t(v)), p_t(v) \rangle$$

- 3 Receives messages from his neighbors, forwarding those that are not older than  $d$



# The communication protocol with delay $d$

For each  $t = 1, \dots, T$  each agent  $v \in V$  does the following:

- 1 Plays an action  $I_t(v)$  drawn according to his private distribution  $p_t(v)$  observing loss  $\ell_t(I_t(v))$  (same loss vector for all agents)
- 2 Sends to his neighbors the message

$$m_t(v) = \langle t, v, I_t(v), \ell_t(I_t(v)), p_t(v) \rangle$$

- 3 Receives messages from his neighbors, forwarding those that are not older than  $d$
- An agent receives a message from another agent with a delay equal to the shortest path between them
  - A message sent by some agent  $v$  at time  $t$  will be received by all agents whose shortest-path distance from  $v$  is at most  $d$



# Average welfare regret

$$R_T^{\text{coop}} = \frac{1}{N} \sum_{v \in V} \mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t(v)) \right] - \min_{i=1, \dots, K} \sum_{t=1}^T \ell_t(i)$$

## Remarks

- Clearly,  $R_T^{\text{coop}} \leq \sqrt{TK \ln K}$  when agents run vanilla Exp3 (no cooperation)
- By using other agent's plays, each agent may estimate  $\ell_t$  better (thus learning nearly at full info rate)
- In general,  $d$  trades off between **quality** and **quantity** of information



# Key inequality

$$\mathbb{E} \left[ \sum_{v \in V} \sum_{i=1}^K \frac{\mathbb{P}_t(I_t(v) = i)}{\mathbb{P}_{t-d}(\ell_{t-d}(i) \text{ is observed by } v)} \right] \leq \frac{e}{1 + e^{-1}} (K \alpha_d + N)$$

$\alpha_d$  is the independence number of the graph obtained from  $G$  by connecting any two vertices whose shortest path distance is at most  $d$



# Average welfare regret bound

$$R_T^{\text{coop}} \leq \sqrt{\underbrace{\left( (d+1) + \frac{K}{N} \alpha_d \right)}_{\text{main term}} \underbrace{T \ln K}_{\text{unavoidable}}}$$



# Average welfare regret bound

$$R_T^{\text{coop}} \leq \sqrt{\underbrace{\left( (d+1) + \frac{K}{N} \alpha_d \right)}_{\text{main term}} \underbrace{T \ln K}_{\text{unavoidable}}}$$

$d = K^{1/2}$  with any connected graph  $G$

- Then  $\alpha_d \leq \frac{2N}{d+2}$  and

$$R_T^{\text{coop}} \leq \sqrt{K^{1/2} T \ln K}$$

- This is better than  $\sqrt{KT}$  (minimax for non-cooperating bandits)



- **Feedback graphs** provide a unified model for studying experts, bandits, and more general settings





# Final remarks

- **Feedback graphs** provide a unified model for studying experts, bandits, and more general settings
- We can characterize the **minimax regret** achievable on any directed feedback graph in terms of its structure



- **Feedback graphs** provide a unified model for studying experts, bandits, and more general settings
- We can characterize the **minimax regret** achievable on any directed feedback graph in terms of its structure
- Besides observability, graphs can be used to define bandit problems with **smooth losses**



- **Feedback graphs** provide a unified model for studying experts, bandits, and more general settings
- We can characterize the **minimax regret** achievable on any directed feedback graph in terms of its structure
- Besides observability, graphs can be used to define bandit problems with **smooth losses**
- **Communication graphs** provide a natural setting for the study of interacting bandits using tools similar to those developed for the analysis of feedback graphs

