

# Bandit Algorithms for Online Optimization

Nicolò Cesa-Bianchi

Università degli Studi di Milano



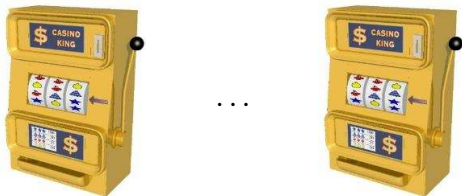
# Multiarmed Bandits



- 1 The multiarmed bandit problem
  - 1 Algorithm Exp3
- 2 Combinatorial bandits (online linear optimization in  $\{0, 1\}^d$ )
  - 1 Spanners and GeometricHedge
  - 2 A more specific approach: ComBand
  - 3 Self-concordant functions
- 3 Discussion



# The multiarmed bandit problem



$N$  slot machines

- Each play of a slot machine (action) returns a **payoff** (gain or loss)
- Design a strategy of **repeated play** to maximize gain / minimize loss
- A classical problem in **sequential design of experiments**  
[Robbins, 1952]
- **Modern applications:** ad placement, dynamic pricing, routing in networks, active model selection, ...



# Main ingredients

- **Partial feedback:** only one action is played at each time step



# Main ingredients

- **Partial feedback:** only one action is played at each time step
- **Exploration:** find out the single best action
- **Exploitation:** play the best action as often as possible



# Main ingredients

- **Partial feedback:** only one action is played at each time step
- **Exploration:** find out the single best action
- **Exploitation:** play the best action as often as possible
- **Dilemma:** when to explore? And how much?



# Main ingredients

- **Partial feedback:** only one action is played at each time step
- **Exploration:** find out the single best action
- **Exploitation:** play the best action as often as possible
- **Dilemma:** when to explore? And how much?

## Payoff generation

- Stochastic/nonstochastic assumptions
- We focus on the **nonstochastic** case





## Problem description

- A fixed but unknown sequence  $x_1, x_2, \dots$  of **loss functions**

$$x_t : \{1, \dots, N\} \rightarrow [0, 1]$$

- $x_t(i)$  is loss of action  $i = 1, \dots, N$  if played at time  $t$



## Problem description

- A fixed but unknown sequence  $x_1, x_2, \dots$  of **loss functions**

$$x_t : \{1, \dots, N\} \rightarrow [0, 1]$$

- $x_t(i)$  is loss of action  $i = 1, \dots, N$  if played at time  $t$

For  $t = 1, 2, \dots$

- 1 Select action  $I_t \in \{1, \dots, N\}$
- 2 Observe loss  $x_t(I_t)$  (losses  $x_t(i)$  for  $i \neq I_t$  remains **hidden**)



## Problem description

- A fixed but unknown sequence  $x_1, x_2, \dots$  of **loss functions**

$$x_t : \{1, \dots, N\} \rightarrow [0, 1]$$

- $x_t(i)$  is loss of action  $i = 1, \dots, N$  if played at time  $t$

For  $t = 1, 2, \dots$

- 1 Select action  $I_t \in \{1, \dots, N\}$
- 2 Observe loss  $x_t(I_t)$  (losses  $x_t(i)$  for  $i \neq I_t$  remains **hidden**)

Regret for play sequence  $I_1, I_2, \dots, I_T$

How much we lose by not consistently playing the **single best** action

$$\sum_{t=1}^T x_t(I_t) - \min_{i=1, \dots, N} \sum_{t=1}^T x_t(i)$$

- **Randomization:** At time  $t$ , action  $i$  is played with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$



- **Randomization:** At time  $t$ , action  $i$  is played with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \underbrace{\sum_{s=1}^{t-1} \hat{x}_s(i)}_{\text{estimated past loss}}\right)$



- **Randomization:** At time  $t$ , action  $i$  is played with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \underbrace{\sum_{s=1}^{t-1} \hat{x}_s(i)}_{\text{estimated past loss}}\right)$

- **Loss vector estimate:**  $\hat{x}_t(i) = \frac{x_t(i)}{p_t(i)} \mathbb{I}_{\{I_t=i\}}$



- **Randomization:** At time  $t$ , action  $i$  is played with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \underbrace{\sum_{s=1}^{t-1} \hat{x}_s(i)}_{\text{estimated past loss}}\right)$

- **Loss vector estimate:**  $\hat{x}_t(i) = \frac{x_t(i)}{p_t(i)} \mathbb{I}_{\{I_t=i\}}$

- **Key properties:**  $\mathbb{E}_t[\hat{x}_t(i)] = x_t(i) \quad |\hat{x}_t(i)| \leq \frac{N}{\gamma}$

- **Randomization:** At time  $t$ , action  $i$  is played with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \underbrace{\sum_{s=1}^{t-1} \hat{x}_s(i)}_{\text{estimated past loss}}\right)$

- **Loss vector estimate:**  $\hat{x}_t(i) = \frac{x_t(i)}{p_t(i)} \mathbb{I}_{\{I_t=i\}}$

- **Key properties:**  $\mathbb{E}_t[\hat{x}_t(i)] = x_t(i)$   $|\hat{x}_t(i)| \leq \frac{N}{\gamma}$



## Theorem

For any sequence  $x_1, x_2, \dots$  of loss functions

$$\sum_{t=1}^T \mathbb{E}[x_t(I_t)] - \min_{i=1, \dots, N} \sum_{t=1}^T x_t(i) \leq c \sqrt{TN \ln N}$$



## Theorem

For any sequence  $x_1, x_2, \dots$  of loss functions

$$\sum_{t=1}^T \mathbb{E}[x_t(I_t)] - \min_{i=1, \dots, N} \sum_{t=1}^T x_t(i) \leq c \sqrt{TN \ln N}$$

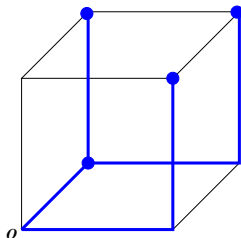
## Remarks

- Factor  $\sqrt{T \ln N}$  is necessary even when  $x_t(i)$  is known for  $i \neq I_t$
- Factor  $\sqrt{N}$  is due to **range of estimates**  $|\hat{x}_t(i)| \leq \frac{N}{\gamma}$



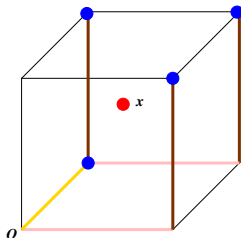
# Combinatorial bandits

- Actions are **boolean vectors**:  $\mathbf{v}(i) \subseteq \{0,1\}^d \quad i = 1, \dots, N$



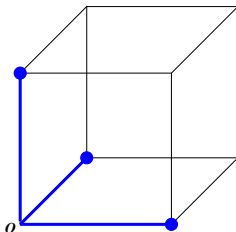
# Combinatorial bandits

- Actions are **boolean vectors**:  $\mathbf{v}(i) \subseteq \{0,1\}^d \quad i = 1, \dots, N$
- Losses  $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^d$  are **linear**:  $\mathbf{v}(i)^\top \mathbf{x}_t$



# Combinatorial bandits

- Actions are **boolean vectors**:  $\mathbf{v}(i) \subseteq \{0,1\}^d \quad i = 1, \dots, N$
- Losses  $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^d$  are **linear**:  $\mathbf{v}(i)^\top \mathbf{x}_t$
- **Classical bandits**:  $\mathbf{v}(1), \dots, \mathbf{v}(N)$  are the canonical basis of  $\mathbb{R}^d$



- Loss vector estimation problem is now  $d$ -dimensional
- Expect regret bound  $\sqrt{Td \ln N}$  instead of  $\sqrt{TN \ln N}$



# Application to multidimensional allocation problems

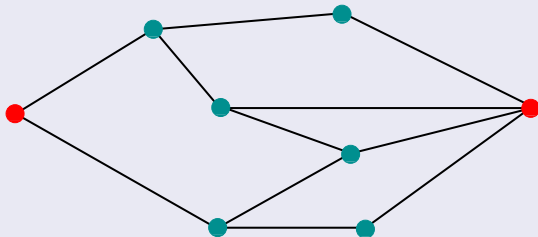
## Multitask bandits

$d = 6$



## Path selection problems on networks

$d = 11$



# Application to multidimensional allocation problems

## Multitask bandits

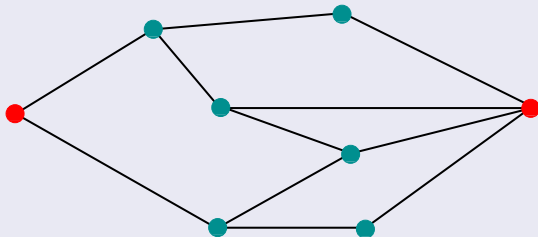
$d = 6$



$$\mathbf{v}(I_t) = (1, 0, 0, 0, 0, 1)$$

## Path selection problems on networks

$d = 11$





# Application to multidimensional allocation problems

## Multitask bandits

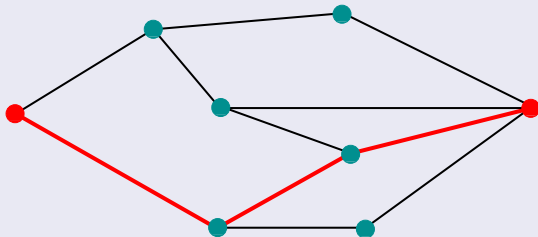
$d = 6$



$$\mathbf{v}(I_t) = (1, 0, 0, 0, 0, 1)$$

## Path selection problems on networks

$d = 11$



# Application to multidimensional allocation problems

## Multitask bandits

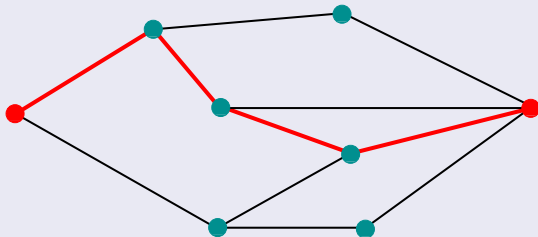
$d = 6$



$$\mathbf{v}(I_t) = (1, 0, 0, 0, 0, 1)$$

## Path selection problems on networks

$d = 11$



# More examples

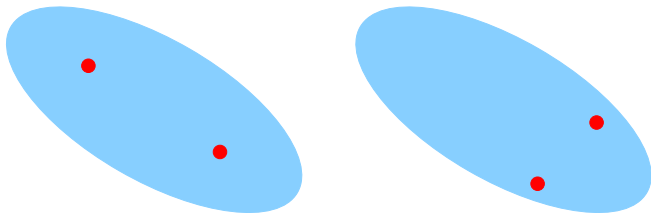
Actions are combinatorial objects represented by their **incidence vectors**  $\mathbf{v}(i) \subseteq \{0, 1\}^d$

$K$ -sized subsets of $d$ elements	$\{0, 1\}^d$	$N = \binom{d}{K}$
Permutations of $K$ objects	$\{0, 1\}^{K^2}$	$N = K!$
Spanning trees of a $K$ -clique	$\{0, 1\}^{\binom{K}{2}}$	$N = K^{K-2}$
Balanced cuts of a $2K$ -clique	$\{0, 1\}^{\binom{2K}{2}}$	$N = \binom{2K}{K}$
Hamiltonian cycles in a $K$ -clique	$\{0, 1\}^{\binom{K}{2}}$	$N = \frac{(K-1)!}{2}$



# Barycentric spanners

- Regret of Exp3 is  $\sqrt{TN \ln N}$  ... bad when  $N$  is large
- Factor  $\sqrt{N}$  comes from **range of estimates**  $N/\gamma$  controlled by exploration term – we are not exploiting overlappings
- **Idea:** explore on a subset of the action set and estimate losses of remaining actions by interpolation



## Theorem

For any  $\mathcal{S} \subseteq \mathbb{R}^d$  there exists a (generally nonorthogonal) basis for the span of  $\mathcal{S}$  such that  $\|\mathbf{v}\|_\infty \leq 1$  for all  $\mathbf{v} \in \mathcal{S}$



## Theorem

For any  $\mathcal{S} \subseteq \mathbb{R}^d$  there exists a (generally nonorthogonal) basis for the span of  $\mathcal{S}$  such that  $\|\mathbf{v}\|_\infty \leq 1$  for all  $\mathbf{v} \in \mathcal{S}$

## Coordinate-free scaling assumption

For any  $\mathbf{v} \in \mathcal{S}$        $\mathbf{v}^\top \mathbf{x}_t \leq B$



- Draw  $\mathbf{v}(I_t) \in \{\mathbf{v}(1), \dots, \mathbf{v}(N)\}$  with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \frac{\gamma}{d} \underbrace{\mathbb{I}_{\{\mathbf{v}(i) \text{ is spanner}\}}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- Draw  $\mathbf{v}(I_t) \in \{\mathbf{v}(1), \dots, \mathbf{v}(N)\}$  with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \frac{\gamma}{d} \underbrace{\mathbb{I}_{\{\mathbf{v}(i) \text{ is spanner}\}}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- Observe loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$



- Draw  $\mathbf{v}(I_t) \in \{\mathbf{v}(1), \dots, \mathbf{v}(N)\}$  with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \frac{\gamma}{d} \underbrace{\mathbb{I}_{\{\mathbf{v}(i) \text{ is spanner}\}}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- Observe loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$
- **Loss vector estimate** (least squares):

$$\hat{\mathbf{x}}_t = \mathbf{C}_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t \quad \text{where} \quad \mathbf{C}_t = \underbrace{\mathbb{E}_t[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top]}_{\text{correlation matrix}}$$

- Draw  $\mathbf{v}(I_t) \in \{\mathbf{v}(1), \dots, \mathbf{v}(N)\}$  with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \frac{\gamma}{d} \underbrace{\mathbb{I}_{\{\mathbf{v}(i) \text{ is spanner}\}}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- Observe loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$

- **Loss vector estimate** (least squares):

$$\hat{\mathbf{x}}_t = \mathbf{C}_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t \quad \text{where} \quad \mathbf{C}_t = \underbrace{\mathbb{E}_t[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top]}_{\text{correlation matrix}}$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \sum_{s=1}^{t-1} \hat{\mathbf{x}}_s^\top \mathbf{v}(i)\right)$

- Loss is unbiased

$$\mathbb{E}_t[\hat{\mathbf{x}}_t] = \mathbb{E}_t\left[C_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t\right] = C_t^+ \mathbb{E}_t\left[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top\right] \mathbf{x}_t = \mathbf{x}_t$$



- Loss is unbiased

$$\mathbb{E}_t[\widehat{\mathbf{x}}_t] = \mathbb{E}_t\left[C_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t\right] = C_t^+ \mathbb{E}_t\left[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top\right] \mathbf{x}_t = \mathbf{x}_t$$

- Size of estimated losses

$$|\widehat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\text{scalar}} |\mathbf{v}(I_t)^\top C_t^+ \mathbf{v}(i)| \leq \underbrace{\|C_t^+\|}_{\text{matrix norm}} \underbrace{\left(\max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2\right)}_{\text{max norm squared}}$$



- Loss is unbiased

$$\mathbb{E}_t[\widehat{\mathbf{x}}_t] = \mathbb{E}_t\left[\mathbf{C}_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t\right] = \mathbf{C}_t^+ \mathbb{E}_t\left[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top\right] \mathbf{x}_t = \mathbf{x}_t$$

- Size of estimated losses

$$|\widehat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\leq B} |\mathbf{v}(I_t)^\top \mathbf{C}_t^+ \mathbf{v}(i)| \leq \underbrace{\|\mathbf{C}_t^+\|}_{\leq B} \underbrace{\left(\max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2\right)}_{\leq B}$$

- Using scaling assumption  $\mathbf{v}^\top \mathbf{x}_t \leq B$



# GEOMETRIC HEDGE — Properties

- Loss is unbiased

$$\mathbb{E}_t[\widehat{\mathbf{x}}_t] = \mathbb{E}_t\left[C_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t\right] = C_t^+ \mathbb{E}_t\left[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top\right] \mathbf{x}_t = \mathbf{x}_t$$

- Size of estimated losses

$$|\widehat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\leq B} |\mathbf{v}(I_t)^\top C_t^+ \mathbf{v}(i)| \leq \underbrace{\|C_t^+\|}_{\leq d/\gamma} \underbrace{\left(\max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2\right)}_{\leq d}$$

- Using scaling assumption  $\mathbf{v}^\top \mathbf{x}_t \leq B$
- Using spanner property  $\|\mathbf{v}\|_\infty \leq 1$



# GEOMETRIC HEDGE — Properties

- Loss is unbiased

$$\mathbb{E}_t[\widehat{\mathbf{x}}_t] = \mathbb{E}_t\left[C_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t\right] = C_t^+ \mathbb{E}_t\left[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top\right] \mathbf{x}_t = \mathbf{x}_t$$

- Size of estimated losses

$$|\widehat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\leq B} |\mathbf{v}(I_t)^\top C_t^+ \mathbf{v}(i)| \leq \underbrace{\|C_t^+\|}_{\leq d/\gamma} \underbrace{\left(\max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2\right)}_{\leq d}$$

- Using scaling assumption  $\mathbf{v}^\top \mathbf{x}_t \leq B$
- Using spanner property  $\|\mathbf{v}\|_\infty \leq 1$

Final regret bound

$$B d \sqrt{T \ln N}$$



- Draw  $\mathbf{v}(I_t) \in \{\mathbf{v}(1), \dots, \mathbf{v}(N)\}$  with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \frac{\gamma}{d} \underbrace{\mathbb{I}_{\{\mathbf{v}(i) \text{ is spanner}\}}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- Observe loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$

- **Loss vector estimate:**

$$\hat{\mathbf{x}}_t = C_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t \quad \text{where} \quad C_t = \mathbb{E}_t[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top]$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \sum_{s=1}^{t-1} \hat{\mathbf{x}}_s^\top \mathbf{v}(i)\right)$





- Draw  $\mathbf{v}(I_t) \in \{\mathbf{v}(1), \dots, \mathbf{v}(N)\}$  with probability

$$p_t(i) = (1 - \gamma) \underbrace{\frac{w_t(i)}{\sum_j w_t(j)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}} \quad (0 \leq \gamma \leq 1)$$

- Observe loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$

- **Loss vector estimate:**

$$\hat{\mathbf{x}}_t = C_t^+ \mathbf{v}(I_t) \mathbf{v}(I_t)^\top \mathbf{x}_t \quad \text{where} \quad C_t = \mathbb{E}_t[\mathbf{v}(I_t) \mathbf{v}(I_t)^\top]$$

- **Weights:**  $w_t(i) = \exp\left(-\frac{\gamma}{N} \sum_{s=1}^{t-1} \hat{\mathbf{x}}_s^\top \mathbf{v}(i)\right)$



## Coordinate-dependent scaling assumption

For  $\mathcal{S} \subseteq \{0, 1\}^d$        $\|\mathbf{v}\|_1 \leq B$        $\|\mathbf{x}_t\|_\infty \leq 1$



# COMBAND — Properties

## Coordinate-dependent scaling assumption

For  $\mathcal{S} \subseteq \{0, 1\}^d$        $\|\mathbf{v}\|_1 \leq B$        $\|\mathbf{x}_t\|_\infty \leq 1$

## Size of estimated losses

$$|\hat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\leq B} \underbrace{|\mathbf{v}(I_t)^\top \mathbf{C}_t^+ \mathbf{v}(i)|}_{\leq \|\mathbf{C}_t^+\|} \underbrace{\left( \max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2 \right)}_{\leq B^2}$$



## Coordinate-dependent scaling assumption

For  $\mathcal{S} \subseteq \{0, 1\}^d$        $\|\mathbf{v}\|_1 \leq B$        $\|\mathbf{x}_t\|_\infty \leq 1$

## Size of estimated losses

$$|\hat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(\mathbf{I}_t)^\top \mathbf{x}_t}_{\leq B} |\mathbf{v}(\mathbf{I}_t)^\top \mathbf{C}_t^+ \mathbf{v}(i)| \leq B \underbrace{\|\mathbf{C}_t^+\|}_{\leq B} \underbrace{\left( \max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2 \right)}$$

- Using scaling assumption  $\|\mathbf{v}\|_1 \leq B$  and  $\|\mathbf{x}_t\|_\infty \leq 1$



## Coordinate-dependent scaling assumption

For  $\mathcal{S} \subseteq \{0, 1\}^d$        $\|\mathbf{v}\|_1 \leq B$        $\|\mathbf{x}_t\|_\infty \leq 1$

## Size of estimated losses

$$|\hat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\leq B} |\mathbf{v}(I_t)^\top \mathbf{C}_t^+ \mathbf{v}(i)| \leq B \underbrace{\|\mathbf{C}_t^+\|}_{\leq 1/(\gamma \lambda_{\min})} \underbrace{\left( \max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2 \right)}$$

- Using scaling assumption  $\|\mathbf{v}\|_1 \leq B$  and  $\|\mathbf{x}_t\|_\infty \leq 1$
- $\lambda_{\min}$  = smallest eigenval of  $\mathbb{E}[\mathbf{v}(I_t)\mathbf{v}(I_t)^\top]$  w.r.t. the **uniform dist.**



## Coordinate-dependent scaling assumption

For  $\mathcal{S} \subseteq \{0, 1\}^d$        $\|\mathbf{v}\|_1 \leq B$        $\|\mathbf{x}_t\|_\infty \leq 1$

## Size of estimated losses

$$|\hat{\mathbf{x}}_t^\top \mathbf{v}(i)| = \underbrace{\mathbf{v}(I_t)^\top \mathbf{x}_t}_{\leq B} |\mathbf{v}(I_t)^\top \mathbf{C}_t^+ \mathbf{v}(i)| \leq B \underbrace{\|\mathbf{C}_t^+\|}_{\leq 1/(\gamma \lambda_{\min})} \underbrace{\left( \max_{\mathbf{v} \in \mathcal{S}} \|\mathbf{v}\|^2 \right)}_{\leq B}$$

- Using scaling assumption  $\|\mathbf{v}\|_1 \leq B$  and  $\|\mathbf{x}_t\|_\infty \leq 1$
- $\lambda_{\min}$  = smallest eigenval of  $\mathbb{E}[\mathbf{v}(I_t)\mathbf{v}(I_t)^\top]$  w.r.t. the **uniform dist.**
- $\mathbf{v} \in \{0, 1\}^d \rightarrow \|\mathbf{v}\|^2 \leq \|\mathbf{v}\|_1 \leq B$



## Regret bound

$$B \sqrt{\left(1 + \frac{B}{d \lambda_{\min}}\right) T d \ln N}$$



## Regret bound

$$B \sqrt{\left(1 + \frac{B}{d \lambda_{\min}}\right) T d \ln N}$$

## Theorem

In all previous examples\* of combinatorial bandits the condition  $\lambda_{\min} = \Omega(B/d)$  holds

---

\* but paths! ☹️





## Regret bound

$$B \sqrt{\left(1 + \frac{B}{d \lambda_{\min}}\right) T d \ln N}$$

## Theorem

In all previous examples\* of combinatorial bandits the condition  $\lambda_{\min} = \Omega(B/d)$  holds

---

\* but paths! 😞

## Corollary

- In these examples the regret becomes  $O\left(B \sqrt{T d \ln N}\right)$
- This is **not improvable** in general for  $\mathcal{S} \subseteq \{0, 1\}^d$

# How many parameters?

## Recall

Loss vector estimates  $\hat{\mathbf{x}}_t = \mathbf{C}_t^+ \mathbf{v}(\mathbf{I}_t) \mathbf{v}(\mathbf{I}_t)^\top \mathbf{x}_t$

- Maintain  $d$  weights over coordinates  $j = 1, \dots, d$

$$w_{j,t} = \exp \left( -\frac{\gamma}{N} \sum_{s=1}^{t-1} \hat{x}_{j,t} \right)$$



# How many parameters?

## Recall

Loss vector estimates  $\hat{\mathbf{x}}_t = \mathbf{C}_t^+ \mathbf{v}(\mathbf{I}_t) \mathbf{v}(\mathbf{I}_t)^\top \mathbf{x}_t$

- Maintain  $d$  weights over coordinates  $j = 1, \dots, d$

$$w_{j,t} = \exp \left( -\frac{\gamma}{N} \sum_{s=1}^{t-1} \hat{x}_{j,t} \right)$$

- ... translate them into weights over actions  $i = 1, \dots, N$

$$w_t(i) = \prod_{j: v_j(i)=1} w_{j,t}$$



# How many parameters?

## Recall

Loss vector estimates  $\hat{\mathbf{x}}_t = \mathbf{C}_t^+ \mathbf{v}(\mathbf{I}_t) \mathbf{v}(\mathbf{I}_t)^\top \mathbf{x}_t$

- Maintain  $d$  weights over coordinates  $j = 1, \dots, d$

$$w_{j,t} = \exp \left( -\frac{\gamma}{N} \sum_{s=1}^{t-1} \hat{x}_{j,t} \right)$$

- ... translate them into weights over actions  $i = 1, \dots, N$

$$w_t(i) = \prod_{j: v_j(i)=1} w_{j,t}$$

## Problem

How to efficiently sample from  $w_t(1), \dots, w_t(N)$  ?

# Efficient sampling

$K$ bandit problems	😊	independent draws from each bandit
$K$ -sized subsets	😊	compute conditionals using dynamic programming
Permutations of $K$ objects	😊	random sampling of perfect matchings
Spanning trees of a $K$ -clique	😞	OK for uniform, unknown for weighted
Balanced cuts of a $2K$ -clique	😊	sampling ferromagnetic Ising model
Hamiltonian cycles in a $K$ -clique	😞	notoriously hard



# A more general framework

## Reduction to online gradient descent

[Zinkevich, 2003]

- 1 Pick  $\mathbf{p}_t$  from the **convex hull** of  $\mathcal{S}$  and use it to draw action  $I_t \in \mathcal{S}$
- 2 Suffer loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$  and observe loss  $\mathbf{p}_t^\top \mathbf{x}_t$
- 3 Compute loss vector estimate  $\hat{\mathbf{x}}_t$
- 4 Perform gradient update:  $\mathbf{p}_{t+1} = \mathbf{p}_t - \eta \hat{\mathbf{x}}_t$



# A more general framework

## Reduction to online gradient descent

[Zinkevich, 2003]

- 1 Pick  $\mathbf{p}_t$  from the **convex hull** of  $\mathcal{S}$  and use it to draw action  $I_t \in \mathcal{S}$
- 2 Suffer loss  $\mathbf{v}(I_t)^\top \mathbf{x}_t$  and observe loss  $\mathbf{p}_t^\top \mathbf{x}_t$
- 3 Compute loss vector estimate  $\hat{\mathbf{x}}_t$
- 4 Perform gradient update:  $\mathbf{p}_{t+1} = \mathbf{p}_t - \eta \hat{\mathbf{x}}_t$

## Issues

**Sampling problem:** need to draw  $I_t$  so that  $\mathbb{E}[\mathbf{v}(I_t)^\top \mathbf{x}_t] = \mathbf{p}_t^\top \mathbf{x}_t$

Point  $\mathbf{p}_t$  must simultaneously

- 1 have small loss  $\mathbf{p}_t^\top \mathbf{x}_t$  (exploitation)
- 2 lead to a good estimate  $\hat{\mathbf{x}}_t$  (exploration)

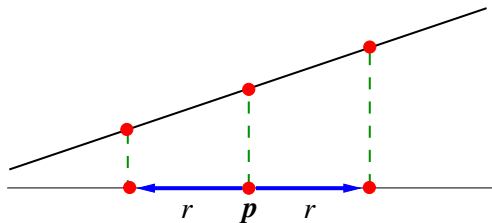


# Gradient descent without a gradient

[Flaxman, Kalai and McMahan, 2004]

- Play  $\mathbf{p}_t + r\mathbf{u}$  where  $\mathbf{u}$  is a random unit vector and  $r > 0$
- Loss vector estimate:  $\hat{\mathbf{x}}_t = d(\mathbf{p}_t + r\mathbf{u})^\top \mathbf{x}_t \frac{\mathbf{u}}{r}$
- In one dimension

$$\mathbb{E}[\hat{\mathbf{x}}_t] = \frac{1}{2} \frac{(\mathbf{p}_t + r)\mathbf{x}_t}{r} + \frac{1}{2} \frac{(\mathbf{p}_t - r)\mathbf{x}_t}{r} = \mathbf{x}_t$$

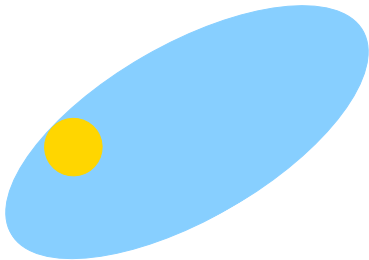




# Issues

- $r \in \mathbb{R}$  in  $\mathbf{p}_t + r\mathbf{u}$  determines the sizes of estimated losses
- $r$  cannot be too big when  $\mathbf{p}_t$  is close to a border
- $r$  must be big enough to control the estimates

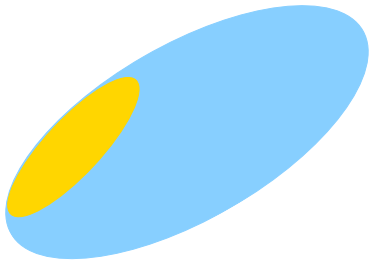
$$\hat{\mathbf{x}}_t = d(\mathbf{p}_t + r\mathbf{u})^\top \mathbf{x}_t \frac{\mathbf{u}}{r}$$



- $r \in \mathbb{R}$  in  $\mathbf{p}_t + r\mathbf{u}$  determines the sizes of estimated losses
- $r$  cannot be too big when  $\mathbf{p}_t$  is close to a border
- $r$  must be big enough to control the estimates

$$\hat{\mathbf{x}}_t = d(\mathbf{p}_t + r\mathbf{u})^\top \mathbf{x}_t \frac{\mathbf{u}}{r}$$

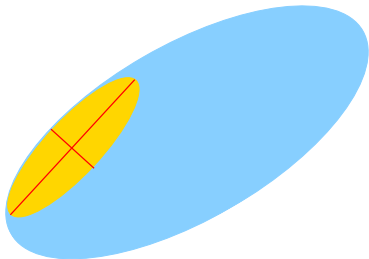
- We need to exploit the **local geometry**



- Barrier functions in interior point optimization
- **Fact:** any convex closed set  $\mathcal{K}$  admits a self-concordant function  $F$  (may be hard to find)



- Barrier functions in interior point optimization
- **Fact:** any convex closed set  $\mathcal{K}$  admits a self-concordant function  $F$  (may be hard to find)
- The Hessian matrix  $\nabla^2 F$  of a self-concordant function  $F$  defines the Dikin ellipsoid which is always contained in  $\mathcal{K}$  (our action space)
- **Loss vector estimate:**  $\hat{\mathbf{x}}_t = \mathbf{p}_t \pm \mathbf{e}_i \sqrt{\lambda_i}$   
where  $\{\mathbf{e}_i, \lambda_i\}$  is a random eigenvector-eigenvalue pair of  $\nabla^2 F$



# Summary for combinatorial bandits

## Regrets

GEOMETRIC HEDGE	COMBAND	SELF CONCORDANT
$Bd \sqrt{T \ln N}$	$B \sqrt{Td \ln N}$	$Bd^{3/2} \sqrt{T \ln T}$



# Summary for combinatorial bandits

## Regrets

GEOMETRIC HEDGE	COMBAND	SELF CONCORDANT
$Bd \sqrt{T \ln N}$	$B \sqrt{Td \ln N}$	$Bd^{3/2} \sqrt{T \ln T}$

## Efficiency

- **All**: efficiency of the associated sampling problem
- **SELF CONCORDANT** efficiency of computation of the self-concordant function



# Summary for combinatorial bandits

## Regrets

GEOMETRIC HEDGE	COMBAND	SELF CONCORDANT
$Bd \sqrt{T \ln N}$	$B \sqrt{Td \ln N}$	$Bd^{3/2} \sqrt{T \ln T}$

## Efficiency

- **All:** efficiency of the associated sampling problem
- **SELF CONCORDANT** efficiency of computation of the self-concordant function

## Remark

Typically:  $\ln N = \Theta(\sqrt{d} \ln d)$



# Summary for combinatorial bandits

## Regrets

GEOMETRICHEDGE	COMBAND	SELFCONCORDANT
$Bd^{5/4} \sqrt{T \ln d}$	$Bd^{3/4} \sqrt{T \ln d}$	$Bd^{3/2} \sqrt{T \ln T}$

## Efficiency

- **All**: efficiency of the associated sampling problem
- **SELFCONCORDANT** efficiency of computation of the self-concordant function

## Remark

Typically:  $\ln N = \mathcal{O}(\sqrt{d} \ln d)$





# Conclusions

- When  $\mathcal{S} \subseteq \{0, 1\}^d$ , COMBAND improves on GEOMETRICHEDGE and SELFCONCORDANT by exploiting structure of  $\mathcal{S}$
- In problems like **path selection**  $\lambda_{\min}$  may get too small  
→ COMBAND is suboptimal
- **Fix in progress:** Perform exploration using a distribution that maximizes  $\lambda_{\min}$
- Is this optimal in general?

