

adversarial bandit problems: the power of randomization

Gábor Lugosi

ICREA and Pompeu Fabra University, Barcelona

on-line prediction

A game between forecaster and environment.

At each round t ,

on-line prediction

A game between forecaster and environment.

At each round t ,

✱ the forecaster chooses an action $I_t \in \{1, \dots, N\}$;

on-line prediction

A game between forecaster and environment.

At each round t ,

- * the forecaster chooses an action $I_t \in \{1, \dots, N\}$;
(actions are often called experts)

on-line prediction

A game between forecaster and environment.

At each round t ,

- * the forecaster chooses an action $I_t \in \{1, \dots, N\}$;
(actions are often called experts)
- * the environment chooses losses $\ell_t(1), \dots, \ell_t(N) \in [0, 1]$;

on-line prediction

A game between forecaster and environment.

At each round t ,

- * the forecaster chooses an action $I_t \in \{1, \dots, N\}$;
(actions are often called experts)
- * the environment chooses losses $\ell_t(1), \dots, \ell_t(N) \in [0, 1]$;
- * the forecaster suffers loss $\ell_t(I_t)$.

on-line prediction

A game between forecaster and environment.

At each round t ,

✱ the forecaster chooses an action $I_t \in \{1, \dots, N\}$;

(actions are often called experts)

✱ the environment chooses losses $\ell_t(1), \dots, \ell_t(N) \in [0, 1]$;

✱ the forecaster suffers loss $\ell_t(I_t)$.

The goal is to minimize the average regret

$$R_n = \frac{1}{n} \left(\sum_{t=1}^n \ell_t(I_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right).$$

outcome sequence

Often $\ell_t(\mathbf{i}) = \ell(\mathbf{i}, \mathbf{y}_t)$

where $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathcal{Y}$ is the **sequence of outcomes** to be predicted.

and $\ell : \{1, \dots, N\} \times \mathcal{Y} \rightarrow [0, 1]$ is a **loss function**.

simplest example

Is it possible to make regret $\rightarrow 0$ for all loss assignments?

simplest example

Is it possible to make regret $\rightarrow 0$ for all loss assignments?

Let $N = 2$ and define, for all $t = 1, \dots, n$,

$$l_t(\mathbf{1}) = \begin{cases} 0 & \text{if } I_t = 2 \\ 1 & \text{if } I_t = 1 \end{cases}$$

and $l_t(\mathbf{2}) = 1 - l_t(\mathbf{1})$.

simplest example

Is it possible to make regret $\rightarrow 0$ for all loss assignments?

Let $N = 2$ and define, for all $t = 1, \dots, n$,

$$l_t(\mathbf{1}) = \begin{cases} 0 & \text{if } I_t = 2 \\ 1 & \text{if } I_t = 1 \end{cases}$$

and $l_t(\mathbf{2}) = 1 - l_t(\mathbf{1})$.

Then

$$\sum_{t=1}^n l_t(I_t) = n \quad \text{and} \quad \min_{i=1,2} \sum_{t=1}^n l_t(i) \leq \frac{n}{2}$$

simplest example

Is it possible to make regret $\rightarrow 0$ for all loss assignments?

Let $N = 2$ and define, for all $t = 1, \dots, n$,

$$l_t(1) = \begin{cases} 0 & \text{if } I_t = 2 \\ 1 & \text{if } I_t = 1 \end{cases}$$

and $l_t(2) = 1 - l_t(1)$.

Then

$$\sum_{t=1}^n l_t(I_t) = n \quad \text{and} \quad \min_{i=1,2} \sum_{t=1}^n l_t(i) \leq \frac{n}{2}$$

so

$$R_n \geq \frac{1}{2} \cdot n$$

randomized prediction

Key to solution: **randomization**.

At time \mathbf{t} , the forecaster chooses a probability distribution

$$\mathbf{p}_{\mathbf{t}-1} = (\mathbf{p}_{1,\mathbf{t}-1}, \dots, \mathbf{p}_{\mathbf{N},\mathbf{t}-1})$$

and chooses action \mathbf{i} with probability $\mathbf{p}_{\mathbf{i},\mathbf{t}-1}$.

Simplest model: all losses $\ell_{\mathbf{s}}(\mathbf{i})$, $\mathbf{i} = 1, \dots, \mathbf{N}$, $\mathbf{s} < \mathbf{t}$, are observed: **full information**.

randomized prediction

This and related models have been studied in

- * game theory: playing repeated games;
- * information theory: gambling and data compression;
- * statistics: sequential decisions;
- * statistical learning theory: on-line learning;

Hannan and Blackwell

Hannan (1957) and Blackwell (1956) showed that the forecaster has a strategy such that

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{I}_t) - \min_{\mathbf{i} \leq N} \sum_{t=1}^n \ell_t(\mathbf{i}) \right) \rightarrow 0$$

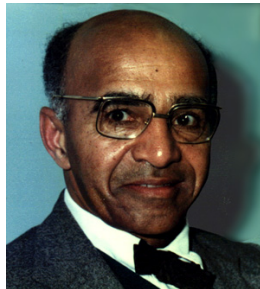
almost surely for all strategies of the environment.

Hannan and Blackwell

Hannan (1957) and Blackwell (1956) showed that the forecaster has a strategy such that

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(I_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) \rightarrow 0$$

almost surely for all strategies of the environment.



basic ideas

expected loss of the forecaster:

$$l_t(\mathbf{p}_{t-1}) = \sum_{i=1}^N p_{i,t-1} l_t(\mathbf{i}) = \mathbb{E}_t l_t(\mathbf{I}_t)$$

basic ideas

expected loss of the forecaster:

$$\ell_t(\mathbf{p}_{t-1}) = \sum_{i=1}^N \mathbf{p}_{i,t-1} \ell_t(\mathbf{i}) = \mathbb{E}_t \ell_t(\mathbf{I}_t)$$

By martingale convergence,

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{I}_t) - \sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) \right) = O_P(n^{-1/2})$$

basic ideas

expected loss of the forecaster:

$$\ell_t(\mathbf{p}_{t-1}) = \sum_{i=1}^N \mathbf{p}_{i,t-1} \ell_t(\mathbf{i}) = \mathbb{E}_t \ell_t(\mathbf{I}_t)$$

By martingale convergence,

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{I}_t) - \sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) \right) = O_P(n^{-1/2})$$

so it suffices to study

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \min_{i \leq N} \sum_{t=1}^n \ell_t(\mathbf{i}) \right)$$

weighted average prediction

Idea: assign a higher probability to better-performing actions.
Vovk (1990), Littlestone and Warmuth (1989).

weighted average prediction

Idea: assign a higher probability to better-performing actions.
Vovk (1990), Littlestone and Warmuth (1989).

A popular choice is

$$p_{i,t-1} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell_s(\mathbf{i})\right)}{\sum_{k=1}^N \exp\left(-\eta \sum_{s=1}^{t-1} \ell_s(\mathbf{k})\right)} \quad \mathbf{i} = 1, \dots, N.$$

where $\eta > 0$.

weighted average prediction

Idea: assign a higher probability to better-performing actions.
Vovk (1990), Littlestone and Warmuth (1989).

A popular choice is

$$p_{i,t-1} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell_s(\mathbf{i})\right)}{\sum_{\mathbf{k}=1}^N \exp\left(-\eta \sum_{s=1}^{t-1} \ell_s(\mathbf{k})\right)} \quad \mathbf{i} = 1, \dots, N.$$

where $\eta > 0$. Then

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \min_{\mathbf{i} \leq N} \sum_{t=1}^n \ell_t(\mathbf{i}) \right) = \sqrt{\frac{\ln N}{2n}}$$

with $\eta = \sqrt{8 \ln N / n}$.

proof

Let $\mathbf{L}_{i,t} = \sum_{s=1}^t \ell_s(\mathbf{i})$ and

$$\mathbf{W}_t = \sum_{i=1}^N \mathbf{w}_{i,t} = \sum_{i=1}^N e^{-\eta \mathbf{L}_{i,t}}$$

for $t \geq 1$, and $\mathbf{W}_0 = \mathbf{N}$.

proof

Let $\mathbf{L}_{i,t} = \sum_{s=1}^t \ell_s(\mathbf{i})$ and

$$\mathbf{W}_t = \sum_{i=1}^N \mathbf{w}_{i,t} = \sum_{i=1}^N e^{-\eta \mathbf{L}_{i,t}}$$

for $t \geq 1$, and $\mathbf{W}_0 = \mathbf{N}$. First observe that

$$\begin{aligned} \ln \frac{\mathbf{W}_n}{\mathbf{W}_0} &= \ln \left(\sum_{i=1}^N e^{-\eta \mathbf{L}_{i,n}} \right) - \ln N \\ &\geq \ln \left(\max_{i=1, \dots, N} e^{-\eta \mathbf{L}_{i,n}} \right) - \ln N \\ &= -\eta \min_{i=1, \dots, N} \mathbf{L}_{i,n} - \ln N . \end{aligned}$$

proof

On the other hand, for each $\mathbf{t} = \mathbf{1}, \dots, \mathbf{n}$

$$\begin{aligned}\ln \frac{W_{\mathbf{t}}}{W_{\mathbf{t}-1}} &= \ln \frac{\sum_{i=1}^N \mathbf{w}_{i,\mathbf{t}-1} e^{-\eta \ell_{\mathbf{t}}(\mathbf{i})}}{\sum_{j=1}^N \mathbf{w}_{j,\mathbf{t}-1}} \\ &\leq -\eta \frac{\sum_{i=1}^N \mathbf{w}_{i,\mathbf{t}-1} \ell_{\mathbf{t}}(\mathbf{i})}{\sum_{j=1}^N \mathbf{w}_{j,\mathbf{t}-1}} + \frac{\eta^2}{8} \\ &= -\eta \ell_{\mathbf{t}}(\mathbf{p}_{\mathbf{t}-1}) + \frac{\eta^2}{8}\end{aligned}$$

by Hoeffding's inequality.

proof

On the other hand, for each $\mathbf{t} = \mathbf{1}, \dots, \mathbf{n}$

$$\begin{aligned}\ln \frac{W_{\mathbf{t}}}{W_{\mathbf{t}-1}} &= \ln \frac{\sum_{i=1}^N \mathbf{w}_{i,\mathbf{t}-1} e^{-\eta \ell_{\mathbf{t}}(\mathbf{i})}}{\sum_{j=1}^N \mathbf{w}_{j,\mathbf{t}-1}} \\ &\leq -\eta \frac{\sum_{i=1}^N \mathbf{w}_{i,\mathbf{t}-1} \ell_{\mathbf{t}}(\mathbf{i})}{\sum_{j=1}^N \mathbf{w}_{j,\mathbf{t}-1}} + \frac{\eta^2}{8} \\ &= -\eta \ell_{\mathbf{t}}(\mathbf{p}_{\mathbf{t}-1}) + \frac{\eta^2}{8}\end{aligned}$$

by Hoeffding's inequality.

Hoeffding (1963): if $\mathbf{X} \in [0, 1]$,

$$\ln \mathbb{E} e^{-\eta \mathbf{X}} \leq -\eta \mathbb{E} \mathbf{X} + \frac{\eta^2}{8}$$

proof

for each $\mathbf{t} = 1, \dots, n$

$$\ln \frac{W_{\mathbf{t}}}{W_{\mathbf{t}-1}} \leq -\eta \ell_{\mathbf{t}}(\mathbf{p}_{\mathbf{t}-1}) + \frac{\eta^2}{8}$$

Summing over $\mathbf{t} = 1, \dots, n$,

$$\ln \frac{W_n}{W_0} \leq -\eta \sum_{\mathbf{t}=1}^n \ell_{\mathbf{t}}(\mathbf{p}_{\mathbf{t}-1}) + \frac{\eta^2}{8} n .$$

Combining these, we get

$$\sum_{\mathbf{t}=1}^n \ell_{\mathbf{t}}(\mathbf{p}_{\mathbf{t}-1}) \leq \min_{i=1, \dots, N} L_{i,n} + \frac{\ln N}{\eta} + \frac{\eta}{8} n$$

lower bound

The upper bound is optimal: for all predictors,

$$\sup_{n, N, \ell_t(i)} \frac{\sum_{t=1}^n \ell_t(I_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i)}{\sqrt{(n/2) \ln N}} \geq 1.$$

lower bound

The upper bound is optimal: for all predictors,

$$\sup_{n, N, \ell_t(\mathbf{i})} \frac{\sum_{t=1}^n \ell_t(\mathbf{I}_t) - \min_{\mathbf{i} \leq N} \sum_{t=1}^n \ell_t(\mathbf{i})}{\sqrt{(n/2) \ln N}} \geq 1.$$

Idea: choose $\ell_t(\mathbf{i})$ to be i.i.d. symmetric Bernoulli coin flips.

lower bound

The upper bound is optimal: for all predictors,

$$\sup_{n, N, \ell_t(i)} \frac{\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i)}{\sqrt{(n/2) \ln N}} \geq 1.$$

Idea: choose $\ell_t(i)$ to be i.i.d. symmetric Bernoulli coin flips.

$$\begin{aligned} & \sup_{\ell_t(i)} \left(\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) \\ & \geq \mathbb{E} \left[\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right] \\ & = \frac{n}{2} - \min_{i \leq N} \mathbf{B}_i \end{aligned}$$

Where $\mathbf{B}_1, \dots, \mathbf{B}_N$ are independent Binomial $(n, 1/2)$.
Use the central limit theorem.

label efficient prediction

The forecaster does not see the outcomes $\ell_t(\mathbf{i})$ unless he asks for them, but can do it only $\mathbf{m} \ll \mathbf{n}$ times.

label efficient prediction

The forecaster does not see the outcomes $\ell_t(\mathbf{i})$ unless he asks for them, but can do it only $m \ll n$ times.

For each round $t = 1, \dots, n$,

- * the environment chooses the losses $\ell_t(\mathbf{i})$ without revealing them;

label efficient prediction

The forecaster does not see the outcomes $\ell_t(\mathbf{i})$ unless he asks for them, but can do it only $m \ll n$ times.

For each round $t = 1, \dots, n$,

- * the environment chooses the losses $\ell_t(\mathbf{i})$ without revealing them;
- * the forecaster chooses \mathbf{p}_{t-1} and draws an action $\mathbf{i}_t \in \{1, \dots, N\}$ according to this distribution;

label efficient prediction

The forecaster does not see the outcomes $\ell_t(\mathbf{i})$ unless he asks for them, but can do it only $m \ll n$ times.

For each round $t = 1, \dots, n$,

- * the environment chooses the losses $\ell_t(\mathbf{i})$ without revealing them;
- * the forecaster chooses \mathbf{p}_{t-1} and draws an action $\mathbf{l}_t \in \{1, \dots, N\}$ according to this distribution;
- * the forecaster incurs loss $\ell_t(\mathbf{l}_t)$ and each action \mathbf{i} incurs loss $\ell_t(\mathbf{i})$. **not of revealed to the forecaster!**;

label efficient prediction

The forecaster does not see the outcomes $\ell_t(\mathbf{i})$ unless he asks for them, but can do it only $\mathbf{m} \ll \mathbf{n}$ times.

For each round $\mathbf{t} = 1, \dots, \mathbf{n}$,

- * the environment chooses the losses $\ell_t(\mathbf{i})$ without revealing them;
- * the forecaster chooses \mathbf{p}_{t-1} and draws an action $\mathbf{l}_t \in \{1, \dots, \mathbf{N}\}$ according to this distribution;
- * the forecaster incurs loss $\ell_t(\mathbf{l}_t)$ and each action \mathbf{i} incurs loss $\ell_t(\mathbf{i})$. *not of revealed to the forecaster!*;
- * the forecaster decides whether he asks for the values of $\ell_t(\mathbf{i})$ if the total number of revealed outcomes up to time $\mathbf{t} - 1$ is less than \mathbf{m} .

a label efficient forecaster

Idea: ask for values randomly (with probability $\approx m/n$) and use the weighted average forecaster with the estimated losses.

a label efficient forecaster

Idea: ask for values randomly (with probability $\approx m/n$) and use the weighted average forecaster with the estimated losses.

Let \mathbf{Z}_t be i.i.d. Bernoulli ϵ ($\approx m/n$).

The forecaster asks for $\ell_t(\mathbf{i})$ iff $\mathbf{Z}_t = \mathbf{1}$.

a label efficient forecaster

Idea: ask for values randomly (with probability $\approx m/n$) and use the weighted average forecaster with the estimated losses.

Let \mathbf{Z}_t be i.i.d. Bernoulli ϵ ($\approx m/n$).

The forecaster asks for $\ell_t(\mathbf{i})$ iff $\mathbf{Z}_t = \mathbf{1}$. Let

$$\tilde{\ell}_t(\mathbf{i}) \stackrel{\text{def}}{=} \begin{cases} \ell_t(\mathbf{i})/\epsilon & \text{if } \mathbf{Z}_t = \mathbf{1}, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

An unbiased estimate!

a label efficient forecaster

Idea: ask for values randomly (with probability $\approx \mathbf{m}/\mathbf{n}$) and use the weighted average forecaster with the estimated losses.

Let \mathbf{Z}_t be i.i.d. Bernoulli ϵ ($\approx \mathbf{m}/\mathbf{n}$).

The forecaster asks for $\ell_t(\mathbf{i})$ iff $\mathbf{Z}_t = \mathbf{1}$. Let

$$\tilde{\ell}_t(\mathbf{i}) \stackrel{\text{def}}{=} \begin{cases} \ell_t(\mathbf{i})/\epsilon & \text{if } \mathbf{Z}_t = \mathbf{1}, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

An unbiased estimate!

For each round $\mathbf{t} = \mathbf{1}, \mathbf{2}, \dots, \mathbf{n}$ draw an action from $\{\mathbf{1}, \dots, \mathbf{N}\}$ according to the distribution

$$\mathbf{p}_{\mathbf{i}, \mathbf{t}-1} = \frac{\exp\left(-\eta \sum_{\mathbf{s}=\mathbf{1}}^{\mathbf{t}-1} \tilde{\ell}_{\mathbf{s}}(\mathbf{i})\right)}{\sum_{\mathbf{k}=\mathbf{1}}^{\mathbf{N}} \exp\left(-\eta \sum_{\mathbf{s}=\mathbf{1}}^{\mathbf{t}-1} \tilde{\ell}_{\mathbf{s}}(\mathbf{k})\right)} \quad \mathbf{i} = \mathbf{1}, \dots, \mathbf{N} .$$

bound for label efficient prediction

With probability at least $1 - \delta$,

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) \leq 9 \sqrt{\frac{\ln N + \ln(4/\delta)}{m}}.$$

(Cesa-Bianchi, Lugosi, Stoltz, 2005)

bound for label efficient prediction

With probability at least $1 - \delta$,

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) \leq 9 \sqrt{\frac{\ln N + \ln(4/\delta)}{m}}.$$

(Cesa-Bianchi, Lugosi, Stoltz, 2005)

Sketch of proof: First bound

$$\sum_{t=1}^n \tilde{\ell}_t(\mathbf{p}_{t-1}) - \min_{i \leq N} \sum_{t=1}^n \tilde{\ell}_t(i)$$

as before.

bound for label efficient prediction

With probability at least $1 - \delta$,

$$\frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) \leq 9 \sqrt{\frac{\ln N + \ln(4/\delta)}{m}}.$$

(Cesa-Bianchi, Lugosi, Stoltz, 2005)

Sketch of proof: First bound

$$\sum_{t=1}^n \tilde{\ell}_t(\mathbf{p}_{t-1}) - \min_{i \leq N} \sum_{t=1}^n \tilde{\ell}_t(i)$$

as before. Then use Bernstein-type martingale inequalities to handle

$$\sum_{t=1}^n \ell_t(\mathbf{l}_t) - \sum_{t=1}^n \tilde{\ell}_t(\mathbf{p}_{t-1})$$

and

$$\min_{i \leq N} \sum_{t=1}^n \tilde{\ell}_t(i) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i)$$

lower bound

For any forecaster asking for at most m values,

$$\sup_{\ell_t(\mathbf{i}) \in \{0,1\}} \mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{I}_t) - \min_{\mathbf{i} \leq N} \sum_{t=1}^n \ell_t(\mathbf{i}) \right) \geq c \sqrt{\frac{\ln N}{m}}.$$

lower bound

For any forecaster asking for at most m values,

$$\sup_{\ell_t(\mathbf{i}) \in \{0,1\}} \mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{I}_t) - \min_{\mathbf{i} \leq \mathbf{N}} \sum_{t=1}^n \ell_t(\mathbf{i}) \right) \geq c \sqrt{\frac{\ln N}{m}}.$$

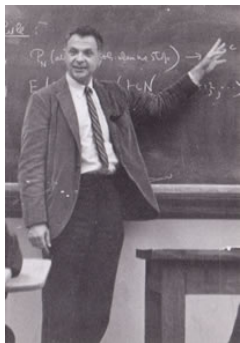
Idea (for $N = 2$): choose the losses randomly (i.i.d.) such that they are either Bernoulli $1/2 - \epsilon$ or Bernoulli $1/2 + \epsilon$.

multi-armed bandits

The forecaster only observes $\ell_t(\mathbf{l}_t)$ but not $\ell_t(\mathbf{i})$ for $\mathbf{i} \neq \mathbf{l}_t$.

multi-armed bandits

The forecaster only observes $l_t(\mathbf{I}_t)$ but not $l_t(\mathbf{i})$ for $\mathbf{i} \neq \mathbf{I}_t$.



Herbert Robbins (1952).

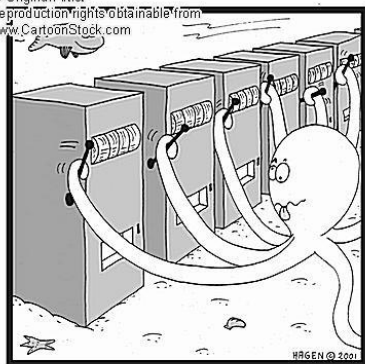
multi-armed bandits



a one-armed bandit

multi-armed bandits

© Original Artist
Reproduction rights obtainable from
www.CartoonStock.com



search ID: cga0170

Compulsive gambling

a multi-armed bandit

multi-armed bandits

Trick: estimate $\ell_t(\mathbf{i})$ by

$$\tilde{\ell}_t(\mathbf{i}) = \frac{\ell_t(\mathbf{l}_t) \mathbb{1}_{\{\mathbf{l}_t = \mathbf{i}\}}}{\mathbf{p}_{\mathbf{l}_t, t-1}}$$

multi-armed bandits

Trick: estimate $\ell_t(\mathbf{i})$ by

$$\tilde{\ell}_t(\mathbf{i}) = \frac{\ell_t(\mathbf{I}_t) \mathbb{1}_{\{\mathbf{I}_t = \mathbf{i}\}}}{\mathbf{p}_{\mathbf{i}, t-1}}$$

This is an unbiased estimate:

$$\mathbb{E}_t \tilde{\ell}_t(\mathbf{i}) = \sum_{j=1}^N \mathbf{p}_{j, t-1} \frac{\ell_t(\mathbf{j}) \mathbb{1}_{\{\mathbf{j} = \mathbf{i}\}}}{\mathbf{p}_{j, t-1}} = \ell_t(\mathbf{i})$$

multi-armed bandits

Trick: estimate $\ell_t(\mathbf{i})$ by

$$\tilde{\ell}_t(\mathbf{i}) = \frac{\ell_t(\mathbf{l}_t) \mathbb{1}_{\{\mathbf{l}_t=\mathbf{i}\}}}{\mathbf{p}_{\mathbf{l}_t, t-1}}$$

This is an unbiased estimate:

$$\mathbb{E}_t \tilde{\ell}_t(\mathbf{i}) = \sum_{\mathbf{j}=1}^N \mathbf{p}_{\mathbf{j}, t-1} \frac{\ell_t(\mathbf{j}) \mathbb{1}_{\{\mathbf{j}=\mathbf{i}\}}}{\mathbf{p}_{\mathbf{j}, t-1}} = \ell_t(\mathbf{i})$$

Use the estimated losses to define exponential weights and mix with uniform (Auer, Cesa-Bianchi, Freund, and Schapire, 2002):

$$\mathbf{p}_{\mathbf{i}, t-1} = (1 - \gamma) \underbrace{\frac{\exp\left(-\eta \sum_{s=1}^{t-1} \tilde{\ell}_s(\mathbf{i})\right)}{\sum_{\mathbf{k}=1}^N \exp\left(-\eta \sum_{s=1}^{t-1} \tilde{\ell}_s(\mathbf{k})\right)}}_{\text{exploitation}} + \underbrace{\frac{\gamma}{N}}_{\text{exploration}}$$

multi-armed bandits

$$\mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) = O \left(\sqrt{\frac{N \ln N}{n}} \right),$$

multi-armed bandits

$$\mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \min_{i \leq N} \sum_{t=1}^n \ell_t(i) \right) = \mathcal{O} \left(\sqrt{\frac{N \ln N}{n}} \right),$$



multi-armed bandits

Lower bound:

$$\sup_{\ell_t(\mathbf{i})} \mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \min_{\mathbf{i} \leq \mathbf{N}} \sum_{t=1}^n \ell_t(\mathbf{i}) \right) \geq \mathbf{C} \sqrt{\frac{\mathbf{N}}{n}},$$

Dependence on \mathbf{N} is not logarithmic anymore!

multi-armed bandits

Lower bound:

$$\sup_{\ell_t(\mathbf{i})} \mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \min_{\mathbf{i} \leq \mathbf{N}} \sum_{t=1}^n \ell_t(\mathbf{i}) \right) \geq \mathbf{C} \sqrt{\frac{\mathbf{N}}{n}},$$

Dependence on \mathbf{N} is not logarithmic anymore!

Audibert and Bubeck (2009) constructed a forecaster with

$$\max_{\mathbf{i} \leq \mathbf{N}} \mathbb{E} \frac{1}{n} \left(\sum_{t=1}^n \ell_t(\mathbf{p}_{t-1}) - \sum_{t=1}^n \ell_t(\mathbf{i}) \right) = \mathbf{O} \left(\sqrt{\frac{\mathbf{N}}{n}} \right),$$

follow the perturbed leader

$$\mathbf{l}_t = \arg \min_{i=1, \dots, N} \sum_{s=1}^{t-1} \ell_s(i) + \mathbf{Z}_{i,t}$$

where the $\mathbf{Z}_{i,t}$ are random noise variables.

By carefully defining the distribution of $\mathbf{Z}_{i,t}$ one can get similar regret bounds for the full information case, [Hannan \(1957\)](#); [Kalai and Vempala \(2003\)](#).

combinatorial experts

Often the class of experts is very large but has some combinatorial structure. **Can the structure be exploited?**

combinatorial experts

Often the class of experts is very large but has some combinatorial structure. Can the structure be exploited?

examples:

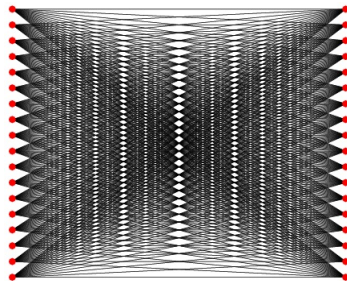
path planning. At each time instance, the forecaster chooses a path in a graph between two fixed nodes. Each edge has an associated loss. Loss of a path is the sum of the losses over the edges in the path.

N is huge!!!



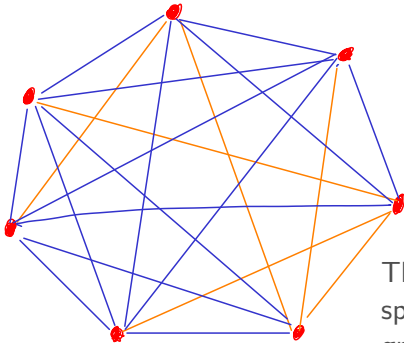
assignments: learning permutations

Given a complete bipartite graph $K_{m,m}$, the forecaster chooses a perfect matching. The loss is the sum of the losses over the edges.



Helmbold and Warmuth (2007): full information case.

spanning trees



The forecaster chooses a spanning tree in the complete graph K_m . The cost is the sum of the losses over the edges.

combinatorial bandits

Two models.

(Easy.) Losses of the components of the chosen object are observed separately. (György, Lugosi, Ottucsák, 2007.)

(Interesting.) Only total loss of the chosen object is observed.

(Awerbuch and Kleinberg, 2004;

McMahan and Blum, 2004;

Dani, Hayes, and Kakade, 2008;

Abernethy, Hazan, and Rakhlin, 2008;

Bartlett, Dani, Hayes, Kakade, and Tewari, 2008;

Cesa-Bianchi and Lugosi, 2009.)

challenges

Performance bounds: Is $O(n^{-1/2}\text{poly}(d))$ regret achievable for the bandit problem?

Algorithmic: How can one draw a random object from the exponentially weighted distribution in polynomial time?

model

$\mathcal{S} = \{\mathbf{v}_1, \dots, \mathbf{v}_N\} \subset \mathbb{R}^d$ is a collection of objects (experts).

Denote $\mathbf{B} = \max_{\mathbf{v}_k} \|\mathbf{v}_k\|_2$.

At every time $\mathbf{t} = 1, 2, \dots$, the opponent chooses a loss vector $\ell_{\mathbf{t}} \in \mathbb{R}^d$.

We assume $\ell_{\mathbf{t}}(\mathbf{k}) = \ell_{\mathbf{t}}^{\top} \mathbf{v}_k \in [-1, 1]$ for all $\mathbf{v}_k \in \mathcal{S}$.

linear bandit problem

For $\mathbf{t} = 1, 2, \dots$,

- * opponent chooses $\ell_t \in \mathbb{R}^d$
- * Forecaster chooses $\mathbf{K}_t \in \{1, \dots, N\}$
- * The cost $\ell_t(\mathbf{K}_t) = \ell_t^\top \mathbf{v}_{\mathbf{K}_t}$ is revealed.

The forecaster's goal is to control the expected **regret**

$$\mathbb{E} \widehat{\mathbf{L}}_n - \min_{\mathbf{k}=1, \dots, N} \mathbf{L}_n(\mathbf{k}) = \sum_{t=1}^n \mathbb{E} \ell_t(\mathbf{K}_t) - \min_{\mathbf{k}=1, \dots, N} \sum_{t=1}^n \ell_t(\mathbf{k}) .$$

Expectation is with respect to the forecaster's internal randomization.

weighted average forecaster

At time \mathbf{t} assign a weight $\mathbf{w}_{\mathbf{t},i}$ to each $\mathbf{i} = 1, \dots, \mathbf{d}$.

The weight of each $\mathbf{v}_{\mathbf{k}} \in \mathcal{S}$ is

$$\bar{\mathbf{w}}_{\mathbf{t}}(\mathbf{k}) = \prod_{\mathbf{i}:\mathbf{v}_{\mathbf{k}}(\mathbf{i})=1} \mathbf{w}_{\mathbf{t},\mathbf{i}} .$$

weighted average forecaster

At time \mathbf{t} assign a weight $\mathbf{w}_{\mathbf{t},i}$ to each $\mathbf{i} = 1, \dots, \mathbf{d}$.

The weight of each $\mathbf{v}_{\mathbf{k}} \in \mathcal{S}$ is

$$\bar{\mathbf{w}}_{\mathbf{t}}(\mathbf{k}) = \prod_{\mathbf{i}: \mathbf{v}_{\mathbf{k}}(\mathbf{i})=1} \mathbf{w}_{\mathbf{t},\mathbf{i}} .$$

Let $\mathbf{q}_{\mathbf{t}-1}(\mathbf{k}) = \bar{\mathbf{w}}_{\mathbf{t}-1}(\mathbf{k}) / \sum_{\mathbf{k}=1}^{\mathbf{N}} \bar{\mathbf{w}}_{\mathbf{t}-1}(\mathbf{k})$.

weighted average forecaster

At time \mathbf{t} assign a weight $\mathbf{w}_{\mathbf{t},i}$ to each $\mathbf{i} = 1, \dots, \mathbf{d}$.

The weight of each $\mathbf{v}_k \in \mathcal{S}$ is

$$\bar{\mathbf{w}}_{\mathbf{t}}(\mathbf{k}) = \prod_{\mathbf{i}: \mathbf{v}_k(\mathbf{i})=1} \mathbf{w}_{\mathbf{t},i}.$$

Let $\mathbf{q}_{\mathbf{t}-1}(\mathbf{k}) = \bar{\mathbf{w}}_{\mathbf{t}-1}(\mathbf{k}) / \sum_{\mathbf{k}=1}^N \bar{\mathbf{w}}_{\mathbf{t}-1}(\mathbf{k})$.

At each time \mathbf{t} , draw $\mathbf{K}_{\mathbf{t}}$ from the distribution

$$\mathbf{p}_{\mathbf{t}-1}(\mathbf{k}) = (1 - \gamma)\mathbf{q}_{\mathbf{t}-1}(\mathbf{k}) + \gamma\boldsymbol{\mu}(\mathbf{k})$$

where $\boldsymbol{\mu}$ is a fixed distribution on \mathcal{S} and $\gamma > 0$. Here

$$\mathbf{w}_{\mathbf{t},i} = \exp(-\eta \tilde{\mathbf{L}}_{\mathbf{t},i})$$

where $\tilde{\mathbf{L}}_{\mathbf{t},i} = \tilde{\ell}_{1,i} + \dots + \tilde{\ell}_{\mathbf{t},i}$ and $\tilde{\ell}_{\mathbf{t},i}$ is an estimated loss.

loss estimates

Dani, Hayes, and Kakade (2008).

Define the scaled incidence vector

$$\mathbf{X}_t = \ell_t(\mathbf{K}_t)\mathbf{V}_{\mathbf{K}_t}$$

where \mathbf{K}_t is distributed according to \mathbf{p}_{t-1} .

loss estimates

Dani, Hayes, and Kakade (2008).

Define the scaled incidence vector

$$\mathbf{X}_t = \ell_t(\mathbf{K}_t) \mathbf{V}_{\mathbf{K}_t}$$

where \mathbf{K}_t is distributed according to \mathbf{p}_{t-1} .

Let $\mathbf{P}_{t-1} = \mathbb{E}[\mathbf{V}_{\mathbf{K}_t} \mathbf{V}_{\mathbf{K}_t}^\top]$ be the $\mathbf{d} \times \mathbf{d}$ correlation matrix.

Hence

$$\mathbf{P}_{t-1}(\mathbf{i}, \mathbf{j}) = \sum_{\mathbf{k} : \mathbf{v}_{\mathbf{k}}(\mathbf{i}) = \mathbf{v}_{\mathbf{k}}(\mathbf{j}) = 1} \mathbf{p}_{t-1}(\mathbf{k}) .$$

Similarly, let \mathbf{Q}_{t-1} and \mathbf{M} be the correlation matrices of $\mathbb{E}[\mathbf{V} \mathbf{V}^\top]$ when \mathbf{V} has law, \mathbf{q}_{t-1} and μ . Then

$$\mathbf{P}_{t-1}(\mathbf{i}, \mathbf{j}) = (1 - \gamma) \mathbf{Q}_{t-1}(\mathbf{i}, \mathbf{j}) + \gamma \mathbf{M}(\mathbf{i}, \mathbf{j}) .$$

loss estimates

Dani, Hayes, and Kakade (2008).

Define the scaled incidence vector

$$\mathbf{X}_t = \ell_t(\mathbf{K}_t)\mathbf{V}_{\mathbf{K}_t}$$

where \mathbf{K}_t is distributed according to \mathbf{p}_{t-1} .

Let $\mathbf{P}_{t-1} = \mathbb{E}[\mathbf{V}_{\mathbf{K}_t} \mathbf{V}_{\mathbf{K}_t}^\top]$ be the $\mathbf{d} \times \mathbf{d}$ correlation matrix.

Hence

$$\mathbf{P}_{t-1}(\mathbf{i}, \mathbf{j}) = \sum_{\mathbf{k} : \mathbf{v}_{\mathbf{k}}(\mathbf{i}) = \mathbf{v}_{\mathbf{k}}(\mathbf{j}) = 1} \mathbf{p}_{t-1}(\mathbf{k}) .$$

Similarly, let \mathbf{Q}_{t-1} and \mathbf{M} be the correlation matrices of $\mathbb{E}[\mathbf{V} \mathbf{V}^\top]$ when \mathbf{V} has law, \mathbf{q}_{t-1} and μ . Then

$$\mathbf{P}_{t-1}(\mathbf{i}, \mathbf{j}) = (1 - \gamma)\mathbf{Q}_{t-1}(\mathbf{i}, \mathbf{j}) + \gamma \mathbf{M}(\mathbf{i}, \mathbf{j}) .$$

The vector of loss estimates is defined by

$$\tilde{\ell}_t = \mathbf{P}_{t-1}^+ \mathbf{X}_t$$

where \mathbf{P}_{t-1}^+ is the pseudo-inverse of \mathbf{P}_{t-1} .

key properties

- * $\mathbf{M} \mathbf{M}^+ \mathbf{v} = \mathbf{v}$ for all $\mathbf{v} \in \mathcal{S}$.
- * \mathbf{Q}_{t-1} is positive semidefinite for every \mathbf{t} .
- * $\mathbf{P}_{t-1} \mathbf{P}_{t-1}^+ \mathbf{v} = \mathbf{v}$ for all \mathbf{t} and $\mathbf{v} \in \mathcal{S}$.

By definition,

$$\mathbb{E}_{\mathbf{t}} \mathbf{X}_{\mathbf{t}} = \mathbf{P}_{t-1} \ell_{\mathbf{t}}$$

and therefore

$$\mathbb{E}_{\mathbf{t}} \tilde{\ell}_{\mathbf{t}} = \mathbf{P}_{t-1}^+ \mathbb{E}_{\mathbf{t}} \mathbf{X}_{\mathbf{t}} = \ell_{\mathbf{t}}$$

An unbiased estimate!

performance bound

The regret of the forecaster satisfies

$$\frac{1}{n} \left(\mathbb{E} \widehat{L}_n - \min_{k=1, \dots, N} L_n(k) \right) \leq 2 \sqrt{\left(\frac{2B^2}{d\lambda_{\min}(\mathbf{M})} + 1 \right) \frac{d \ln N}{n}}.$$

where

$$\lambda_{\min}(\mathbf{M}) = \min_{\mathbf{x} \in \text{span}(\mathcal{S}): \|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{M} \mathbf{x} > 0$$

is the smallest “relevant” eigenvalue of \mathbf{M} . (Cesa-Bianchi and Lugosi, 2009.)

Large $\lambda_{\min}(\mathbf{M})$ is needed to make sure no $|\tilde{\ell}_{t,i}|$ is too large.

performance bound

Other bounds:

$\mathbf{B}\sqrt{\mathbf{d} \ln \mathbf{N}/n}$ (Dani, Hayes, and Kakade). No condition on \mathcal{S} .
Sampling is over a **barycentric spanner**.

$\mathbf{d}\sqrt{(\theta \ln n)/n}$ (Abernethy, Hazan, and Rakhlin). Computationally efficient.

eigenvalue bounds

$$\lambda_{\min}(\mathbf{M}) = \min_{\mathbf{x} \in \text{span}(\mathcal{S}): \|\mathbf{x}\|=1} \mathbb{E}(\mathbf{V}, \mathbf{x})^2 .$$

where \mathbf{V} has distribution μ over \mathcal{S} .

In many cases it suffices to take μ uniform.

multitask bandit problem

The decision maker acts in m games in parallel.

In each game, the decision maker selects one of R possible actions.

After selecting the m actions, the sum of the losses is observed.

$$\lambda_{\min} = \frac{1}{R}$$

$$\max_{\mathbf{k}} \mathbb{E} \left[\widehat{L}_n - L_n(\mathbf{k}) \right] \leq 2m\sqrt{3nR \ln R} .$$

The price of only observing the sum of losses is a factor of m .

Generating a random joint action can be done in polynomial time.

assignments

Perfect matchings of $\mathbf{K}_{m,m}$.

At each time one of the $\mathbf{N} = m!$ perfect matchings of $\mathbf{K}_{m,m}$ is selected.

$$\lambda_{\min}(\mathbf{M}) = \frac{1}{m-1}$$

$$\max_k \mathbb{E} \left[\hat{L}_n - L_n(\mathbf{k}) \right] \leq 2m \sqrt{3n \ln(m!)} .$$

Only a factor of m worse than naive full-information bound.

Sum of weights (partition function) is the permanent of a non-negative matrix. Sampling may be done by a FPAS of Jerrum, Sinclair, and Vigoda (2004).

spanning trees

In a network of m nodes, the cost of communication between two nodes joined by edge e is $\ell_t(e)$ at time t . At each time a minimal connected subnetwork (a spanning tree) is selected. The goal is to minimize the total cost. $N = m^{m-2}$.

$$\lambda_{\min}(\mathbf{M}) = \frac{1}{m} - O\left(\frac{1}{m^2}\right).$$

The entries of \mathbf{M} are

$$\begin{aligned}\mathbb{P}\{\mathbf{V}_i = 1\} &= \frac{2}{m} \\ \mathbb{P}\{\mathbf{V}_i = 1, \mathbf{V}_j = 1\} &= \frac{3}{m^2} \quad \text{if } i \sim j \\ \mathbb{P}\{\mathbf{V}_i = 1, \mathbf{V}_j = 1\} &= \frac{4}{m^2} \quad \text{if } i \not\sim j.\end{aligned}$$

spanning trees

Propp and Wilson (1998) define an exact sampling algorithm. Expected running time is the average hitting time of the Markov chain defined by the edge weights $\mathbf{w}_t(\mathbf{e}) = \exp(-\eta \tilde{L}_t(\mathbf{e}))$.

stars

At each time a central node of a network of m nodes is selected.
Cost is the total cost of the edges adjacent to the node.



$$\lambda_{\min} \geq 1 - O\left(\frac{1}{m}\right).$$

cut sets

A balanced cut in K_{2m} is the collection of all edges between a set of m vertices and its complement. Each balanced cut has m^2 edges and there are $N = \binom{2m}{m}$ balanced cuts.

$$\lambda_{\min}(\mathbf{M}) = \frac{1}{4} - \mathcal{O}\left(\frac{1}{m^2}\right).$$

Choosing from the exponentially weighted average distribution is equivalent to sampling from ferromagnetic Ising model. FPAS by [Randall and Wilson \(1999\)](#).

hamiltonian cycles

A Hamiltonian cycle in K_m is a cycle that visits each vertex exactly once and returns to the starting vertex. $N = (m - 1)!$

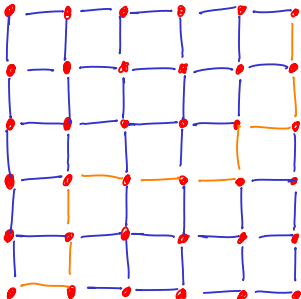
$$\lambda_{\min} \geq \frac{2}{m}$$

Efficient computation is hopeless.

sampling paths

In all these examples μ is uniform over \mathcal{S} .

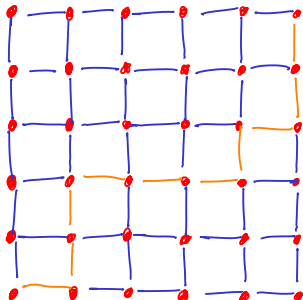
For **path planning** it does not always work.



sampling paths

In all these examples μ is uniform over \mathcal{S} .

For **path planning** it does not always work.



What is the optimal choice of μ ?

What is the optimal way of exploration?

prediction with partial monitoring

For each round $\mathbf{t} = 1, \dots, \mathbf{n}$,

- * the environment chooses the next outcome $\mathbf{J}_t \in \{1, \dots, \mathbf{M}\}$ without revealing it;
- * the forecaster chooses a probability distribution \mathbf{p}_t and draws an action $\mathbf{I}_t \in \{1, \dots, \mathbf{N}\}$ according to \mathbf{p}_t ;
- * the forecaster incurs loss $\ell(\mathbf{I}_t, \mathbf{J}_t)$ and each action \mathbf{i} incurs loss $\ell(\mathbf{i}, \mathbf{J}_t)$. None of these values is revealed to the forecaster;
- * the feedback $\mathbf{h}(\mathbf{I}_t, \mathbf{J}_t)$ is revealed to the forecaster.

$\mathbf{H} = [\mathbf{h}(\mathbf{i}, \mathbf{j})]_{\mathbf{N} \times \mathbf{M}}$ is the feedback matrix.

$\mathbf{L} = [\ell(\mathbf{i}, \mathbf{j})]_{\mathbf{N} \times \mathbf{M}}$ is the loss matrix.

examples

Dynamic pricing. Here $\mathbf{M} = \mathbf{N}$, and $\mathbf{L} = [\ell(\mathbf{i}, \mathbf{j})]_{\mathbf{N} \times \mathbf{N}}$ where

$$\ell(\mathbf{i}, \mathbf{j}) = \frac{(\mathbf{j} - \mathbf{i})\mathbb{1}_{\{\mathbf{i} \leq \mathbf{j}\}} + \mathbf{c}\mathbb{1}_{\{\mathbf{i} > \mathbf{j}\}}}{\mathbf{N}} .$$

and $\mathbf{h}(\mathbf{i}, \mathbf{j}) = \mathbb{1}_{\{\mathbf{i} > \mathbf{j}\}}$ or

$$\mathbf{h}(\mathbf{i}, \mathbf{j}) = \mathbf{a}\mathbb{1}_{\{\mathbf{i} \leq \mathbf{j}\}} + \mathbf{b}\mathbb{1}_{\{\mathbf{i} > \mathbf{j}\}} , \quad \mathbf{i}, \mathbf{j} = 1, \dots, \mathbf{N} .$$

examples

Dynamic pricing. Here $\mathbf{M} = \mathbf{N}$, and $\mathbf{L} = [\ell(\mathbf{i}, \mathbf{j})]_{\mathbf{N} \times \mathbf{N}}$ where

$$\ell(\mathbf{i}, \mathbf{j}) = \frac{(\mathbf{j} - \mathbf{i})\mathbb{1}_{\{\mathbf{i} \leq \mathbf{j}\}} + \mathbf{c}\mathbb{1}_{\{\mathbf{i} > \mathbf{j}\}}}{\mathbf{N}} .$$

and $\mathbf{h}(\mathbf{i}, \mathbf{j}) = \mathbb{1}_{\{\mathbf{i} > \mathbf{j}\}}$ or

$$\mathbf{h}(\mathbf{i}, \mathbf{j}) = \mathbf{a}\mathbb{1}_{\{\mathbf{i} \leq \mathbf{j}\}} + \mathbf{b}\mathbb{1}_{\{\mathbf{i} > \mathbf{j}\}} , \quad \mathbf{i}, \mathbf{j} = 1, \dots, \mathbf{N} .$$

Multi-armed bandit problem. The only information the forecaster receives is his own loss: $\mathbf{H} = \mathbf{L}$.

examples

Apple tasting. $\mathbf{N} = \mathbf{M} = 2$.

$$\mathbf{L} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\mathbf{H} = \begin{bmatrix} a & a \\ b & c \end{bmatrix} .$$

The predictor only receives feedback when he chooses the second action.

examples

Apple tasting. $N = M = 2$.

$$L = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$H = \begin{bmatrix} a & a \\ b & c \end{bmatrix} .$$

The predictor only receives feedback when he chooses the second action.

Label efficient prediction. $N = 3, M = 2$.

$$L = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$H = \begin{bmatrix} a & b \\ c & c \\ c & c \end{bmatrix} .$$

a general predictor

A forecaster first proposed by Piccolboni and Schindelhauer (2001).

Crucial assumption: \mathbf{H} can be encoded such that there exists an $\mathbf{N} \times \mathbf{N}$ matrix $\mathbf{K} = [\mathbf{k}(i, j)]_{\mathbf{N} \times \mathbf{N}}$ such that

$$\mathbf{L} = \mathbf{K} \cdot \mathbf{H} .$$

Thus,

$$\ell(i, j) = \sum_{l=1}^{\mathbf{N}} \mathbf{k}(i, l) \mathbf{h}(l, j) .$$

Then we may estimate the losses by

$$\tilde{\ell}(i, \mathbf{J}_t) = \frac{\mathbf{k}(i, l_t) \mathbf{h}(l_t, \mathbf{J}_t)}{\mathbf{p}_{l_t, t}} .$$

a general predictor

Observe

$$\begin{aligned}\mathbb{E}_t \tilde{\ell}(\mathbf{i}, \mathbf{J}_t) &= \sum_{k=1}^N p_{k,t} \frac{\mathbf{k}(\mathbf{i}, k)h(k, \mathbf{J}_t)}{p_{k,t}} \\ &= \sum_{k=1}^N \mathbf{k}(\mathbf{i}, k)h(k, \mathbf{J}_t) = \ell(\mathbf{i}, \mathbf{J}_t) ,\end{aligned}$$

$\tilde{\ell}(\mathbf{i}, \mathbf{J}_t)$ is an unbiased estimate of $\ell(\mathbf{i}, \mathbf{J}_t)$.

Let

$$p_{\mathbf{i},t} = (1 - \gamma) \frac{e^{-\eta \tilde{\mathbf{L}}_{\mathbf{i},t-1}}}{\sum_{k=1}^N e^{-\eta \tilde{\mathbf{L}}_{k,t-1}}} + \frac{\gamma}{N}$$

where $\tilde{\mathbf{L}}_{\mathbf{i},t} = \sum_{s=1}^t \tilde{\ell}(\mathbf{i}, \mathbf{J}_s)$.

performance bound

With probability at least $1 - \delta$,

$$\begin{aligned} \frac{1}{n} \sum_{t=1}^n \ell(\mathbf{I}_t, \mathbf{J}_t) - \min_{i=1, \dots, N} \frac{1}{n} \sum_{t=1}^n \ell(i, \mathbf{J}_t) \\ \leq \mathbf{C} n^{-1/3} N^{2/3} \sqrt{\ln(N/\delta)}. \end{aligned}$$

where \mathbf{C} depends on \mathbf{K} . (Cesa-Bianchi, Lugosi, Stoltz (2006))

performance bound

With probability at least $1 - \delta$,

$$\begin{aligned} \frac{1}{n} \sum_{t=1}^n \ell(I_t, J_t) - \min_{i=1, \dots, N} \frac{1}{n} \sum_{t=1}^n \ell(i, J_t) \\ \leq C n^{-1/3} N^{2/3} \sqrt{\ln(N/\delta)}. \end{aligned}$$

where C depends on K . (Cesa-Bianchi, Lugosi, Stoltz (2006))

Hannan consistency is achieved with rate $O(n^{-1/3})$ whenever $L = K \cdot H$.

This solves the dynamic pricing problem.

performance bound

With probability at least $1 - \delta$,

$$\begin{aligned} \frac{1}{n} \sum_{t=1}^n \ell(I_t, J_t) - \min_{i=1, \dots, N} \frac{1}{n} \sum_{t=1}^n \ell(i, J_t) \\ \leq C n^{-1/3} N^{2/3} \sqrt{\ln(N/\delta)}. \end{aligned}$$

where C depends on K . (Cesa-Bianchi, Lugosi, Stoltz (2006))

Hannan consistency is achieved with rate $O(n^{-1/3})$ whenever $L = K \cdot H$.

This solves the dynamic pricing problem.

Bartók, Pál, and Szepesvári (2010): if $M = 2$, only possible rates are $n^{-1/2}, n^{-1/3}, 1$

imperfect monitoring: a general framework

\mathbf{S} is a finite set of signals.

Feedback matrix: $\mathbf{H} : \{1, \dots, N\} \times \{1, \dots, M\} \rightarrow \mathcal{P}(\mathbf{S})$.

For each round $\mathbf{t} = 1, 2, \dots, n$,

- * the environment chooses the next outcome $\mathbf{J}_t \in \{1, \dots, M\}$ without revealing it;
- * the forecaster chooses \mathbf{p}_t and draws an action $\mathbf{I}_t \in \{1, \dots, N\}$ according to it;
- * the forecaster receives loss $\ell(\mathbf{I}_t, \mathbf{J}_t)$ and each action \mathbf{i} suffers loss $\ell(\mathbf{i}, \mathbf{J}_t)$, none of these values is revealed to the forecaster;
- * a feedback \mathbf{s}_t drawn at random according to $\mathbf{H}(\mathbf{I}_t, \mathbf{J}_t)$ is revealed to the forecaster.

target

Define

$$\ell(\mathbf{p}, \mathbf{q}) = \sum_{i,j} p_i q_j \ell(i, j)$$

$$\mathbf{H}(\cdot, \mathbf{q}) = (\mathbf{H}(1, \mathbf{q}), \dots, \mathbf{H}(N, \mathbf{q}))$$

where $\mathbf{H}(i, \mathbf{q}) = \sum_j q_j \mathbf{H}(i, j)$.

Denote by \mathcal{F} the set of those Δ that can be written as $\mathbf{H}(\cdot, \mathbf{q})$ for some \mathbf{q} .

\mathcal{F} is the set of “observable” vectors of signal distributions Δ .

The key quantity is

$$\rho(\mathbf{p}, \Delta) = \max_{\mathbf{q}: \mathbf{H}(\cdot, \mathbf{q}) = \Delta} \ell(\mathbf{p}, \mathbf{q})$$

ρ is convex in \mathbf{p} and concave in Δ .

rustichini's theorem

The value of the base one-shot game is

$$\min_{\mathbf{p}} \max_{\mathbf{q}} \ell(\mathbf{p}, \mathbf{q}) = \min_{\mathbf{p}} \max_{\Delta \in \mathcal{F}} \rho(\mathbf{p}, \Delta)$$

If $\bar{\mathbf{q}}_n$ is the empirical distribution of $\mathbf{J}_1, \dots, \mathbf{J}_n$, even with the knowledge of $\mathbf{H}(\cdot, \bar{\mathbf{q}}_n)$ we cannot hope to do better than $\min_{\mathbf{p}} \rho(\mathbf{p}, \mathbf{H}(\cdot, \bar{\mathbf{q}}_n))$.

Rustichini (1999) proved that there exists a strategy such that for all strategies of the opponent, almost surely,

$$\limsup_{n \rightarrow \infty} \left(\frac{1}{n} \sum_{t=1, \dots, n} \ell(\mathbf{I}_t, \mathbf{J}_t) - \min_{\mathbf{p}} \rho(\mathbf{p}, \mathbf{H}(\cdot, \bar{\mathbf{q}}_n)) \right) \leq 0$$

rustichini's theorem

Rustichini's proof relies on an approachability theorem for a continuum of types ([Mertens, Sorin, and Zamir, 1994](#)).

It is non-constructive.

It does not imply any convergence rate.

[Lugosi, Mannor, and Stoltz \(2008\)](#) construct efficiently computable strategies that guarantee fast rates of convergence.