

Strategies for prediction under imperfect monitoring

Gábor Lugosi

ICREA and Pompeu Fabra University, Barcelona, Spain

joint work with

Shie Mannor (Technion, McGill)

Gilles Stoltz (ENS, Paris),

randomized prediction

Prediction game:

number \mathbf{N} of actions,

cardinality \mathbf{M} of outcome space,

reward function $\mathbf{r} : \{1, \dots, \mathbf{N}\} \times \{1, \dots, \mathbf{M}\} \rightarrow [0, 1]$,

number \mathbf{n} of game rounds.

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

- (1) the environment chooses the next outcome \mathbf{J}_t ;
- (2) the forecaster chooses \mathbf{p}_t and determines the random action \mathbf{l}_t , distributed according to \mathbf{p}_t ;
- (3) the environment reveals \mathbf{J}_t ;
- (4) the forecaster receives a reward $\mathbf{r}(\mathbf{l}_t, \mathbf{y}_t)$.

regret

The goal of the forecaster is to minimize the regret

$$\max_{i=1,\dots,N} \frac{1}{n} \sum_{t=1}^n r(i, J_t) - \frac{1}{n} \sum_{t=1}^n r(I_t, J_t)$$

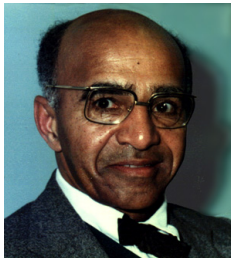
Hannan (1957) and Blackwell (1956) showed that the forecaster has a strategy such that the regret goes to zero almost surely for all strategies of the environment.

regret

The goal of the forecaster is to minimize the regret

$$\max_{i=1,\dots,N} \frac{1}{n} \sum_{t=1}^n r(i, J_t) - \frac{1}{n} \sum_{t=1}^n r(I_t, J_t)$$

Hannan (1957) and Blackwell (1956) showed that the forecaster has a strategy such that the regret goes to zero almost surely for all strategies of the environment.



basic ideas

If $r(\mathbf{p}, \mathbf{j}) = \sum_{i=1}^N p_i r(i, \mathbf{j})$ then by martingale convergence, with probability at least $1 - \delta$,

$$\frac{1}{n} \sum_{t=1}^n r(\mathbf{I}_t, \mathbf{J}_t) \geq \frac{1}{n} \sum_{t=1}^n r(\mathbf{p}_t, \mathbf{J}_t) - \sqrt{\frac{1}{2n} \ln \frac{1}{\delta}}.$$

It suffices to study $(1/n) \sum_{t=1}^n r(\mathbf{p}_t, \mathbf{J}_t)$.

exponential weights

The **exponentially weighted average** forecaster selects action \mathbf{i}_t with probability

$$p_{i,t} = \frac{\exp\left(\eta \sum_{s=1}^{t-1} r(i, \mathbf{J}_s)\right)}{\sum_{k=1}^N \exp\left(\eta \sum_{s=1}^{t-1} r(k, \mathbf{J}_s)\right)}$$

where $\eta > 0$. It is well known that

$$\max_{i=1,\dots,N} \frac{1}{n} \sum_{t=1}^n r(i, \mathbf{J}_t) - \frac{1}{n} \sum_{t=1}^n r(\mathbf{p}_t, \mathbf{J}_t) \leq \frac{\ln N}{n\eta} + \frac{\eta}{8}.$$

With the choice $\eta = \sqrt{8 \ln N / n}$ the upper bound becomes $\sqrt{\ln N / (2n)}$.

Multi-armed bandits

The forecaster only observes $r(l_t, J_t)$ but not $r(i, J_t)$ for $i \neq l_t$.

Trick: estimate $r(i, J_t)$ by

$$\tilde{r}(i, J_t) = \frac{r(l_t, J_t) \mathbb{I}_{[l_t=i]}}{p_{l_t,t}}$$

This is an **unbiased** estimate:

$$\mathbb{E}_t \tilde{r}(i, J_t) = r(i, J_t)$$

Use the estimated losses to define exponential weights. One gets

$$\mathbb{E} \frac{1}{n} \left(\max_{i \leq N} \sum_{t=1}^n r(i, J_t) - \sum_{t=1}^n r(p_t, J_t) \right) = O \left(\sqrt{\frac{N \ln N}{n}} \right),$$

Auer, Cesa-Bianchi, Freund, and Schapire (2002).



imperfect monitoring

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

- 1 the environment chooses the next outcome $\mathbf{J}_t \in \{1, \dots, M\}$ without revealing it;
- 2 the forecaster chooses \mathbf{p}_t and draws an action $\mathbf{I}_t \in \{1, \dots, N\}$ according to it;
- 3 the forecaster receives reward $\mathbf{r}(\mathbf{I}_t, \mathbf{J}_t)$ and each action \mathbf{i} gets reward $\mathbf{r}(\mathbf{i}, \mathbf{J}_t)$, none of these values is revealed to the forecaster;
- 4 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

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

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

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

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

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

rustichini's theorem

The value of the base one-shot game is

$$\max_p \min_q r(p, q) = \max_p \min_{\Delta \in \mathcal{F}} \rho(p, \Delta)$$

If \bar{q}_n is the empirical distribution of y_1, \dots, y_n , even with the knowledge of $H(\cdot, \bar{q}_n)$ we cannot hope to do better than $\max_p \rho(p, H(\cdot, \bar{q}_n))$.

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

$$\limsup_{n \rightarrow \infty} \left(\max_p \rho(p, H(\cdot, \bar{q}_n)) - \frac{1}{n} \sum_{t=1, \dots, n} r(I_t, y_t) \right) \leq 0$$

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.

We construct efficiently computable strategies that guarantee fast rates of convergence.

deterministic feedback, depending on the outcome

Consider first the special case $\mathbf{H}(\mathbf{l}_t, \mathbf{J}_t) = \mathbf{h}(\mathbf{J}_t)$, deterministic. Clearly,

$$\mathbf{r}(\mathbf{p}, \mathbf{j}) \geq \rho(\mathbf{p}, \delta_{\mathbf{h}(\mathbf{j})}) .$$

We introduce a forecaster motivated by the gradient-based strategies (see Cesa-Bianchi and Lugosi, 2006, Section 2.5).

The forecaster uses a sub-gradient of $\rho(\cdot, \delta_{\mathbf{h}(\mathbf{J}_t)})$. $\mathbf{l}_t = \mathbf{i}$ with probability

$$\mathbf{p}_{\mathbf{i}, t} = \frac{e^{\eta \sum_{s=1}^{t-1} (\tilde{\mathbf{r}}(\mathbf{p}_s, \delta_{\mathbf{h}(\mathbf{J}_s))})_{\mathbf{i}}}}{\sum_{\mathbf{j}=1}^{\mathbf{N}} e^{\eta \sum_{s=1}^{t-1} (\tilde{\mathbf{r}}(\mathbf{p}_s, \delta_{\mathbf{h}(\mathbf{J}_s))})_{\mathbf{j}}}} ,$$

where $(\tilde{\mathbf{r}}(\mathbf{p}_s, \delta_{\mathbf{h}(\mathbf{J}_s))})_{\mathbf{i}}$ is the \mathbf{i} -th component of a sub-gradient $\tilde{\mathbf{r}}(\mathbf{p}_s, \delta_{\mathbf{h}(\mathbf{J}_s)}) \in \nabla \rho(\mathbf{p}_s, \delta_{\mathbf{h}(\mathbf{J}_s)})$ of the concave function $\rho(\cdot, \delta_{\mathbf{h}(\mathbf{J}_s)})$.

analysis

$$\begin{aligned} n\rho(\mathbf{p}, \mathbf{H}(\bar{\mathbf{q}}_n)) - \sum_{t=1}^n r(\mathbf{l}_t, \mathbf{J}_t) \\ \approx n\rho(\mathbf{p}, \mathbf{H}(\bar{\mathbf{q}}_n)) - \sum_{t=1}^n r(\mathbf{p}_t, \mathbf{J}_t) \\ \leq n\rho(\mathbf{p}, \mathbf{H}(\bar{\mathbf{q}}_n)) - \sum_{t=1}^n \rho(\mathbf{p}_t, \delta_{\mathbf{h}(\mathbf{J}_t)}) \\ \leq \sum_{t=1}^n (\rho(\mathbf{p}, \delta_{\mathbf{h}(\mathbf{J}_t)}) - \rho(\mathbf{p}_t, \delta_{\mathbf{h}(\mathbf{J}_t)})) \\ \text{(by convexity of } \rho \text{ in the second argument)} \end{aligned}$$

analysis

$$\begin{aligned} &\leq \sum_{t=1}^n \tilde{\mathbf{r}}(\mathbf{p}_t, \delta_{h(J_t)}) \cdot (\mathbf{p} - \mathbf{p}_t) \\ &\quad \text{(by concavity of } \rho \text{ in the first argument)} \\ &\leq \frac{\ln N}{\eta} + \frac{nK^2\eta}{2} \\ &= O(\sqrt{n \ln N}) \end{aligned}$$

random feedback, depending on the outcome

We still assume $\mathbf{H}(\mathbf{I}_t, \mathbf{J}_t) = \mathbf{H}(\mathbf{J}_t)$, but $\mathbf{H}(\mathbf{j})$ is a distribution over signals.

At time t the forecaster observes \mathbf{s}_t drawn from $\mathbf{H}(\mathbf{J}_t)$.

$$\mathbf{H}(\bar{\mathbf{q}}_n) = \frac{1}{n} \sum_{t=1}^n \mathbf{H}(\mathbf{J}_t)$$

needs to be estimated.

Idea: a **lazy** strategy. Group together $m \ll n$ time rounds and use the same mixed strategy.

In the \mathbf{b} -th group one can calculate

$$\frac{1}{m} \sum_{t=bm+1}^{(b+1)m} \delta_{st}$$

and project it to the set \mathcal{F} of feasible distributions:

$$\hat{\Delta}^b = \Pi \left(\frac{1}{m} \sum_{t=bm+1}^{(b+1)m} \delta_{st} \right) .$$

If

$$\Delta^b = \frac{1}{m} \sum_{t=bm+1}^{(b+1)m} H(J_t) \quad (\in \mathcal{F})$$

then $\hat{\Delta}^b \approx \Delta^b$ by (vector valued) martingale convergence.

the strategy

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

- 1 If $\mathbf{b}m + 1 \leq \mathbf{t} < (\mathbf{b} + 1)m$ for some integer \mathbf{b} , choose the distribution $\mathbf{p}_{\mathbf{t}} = \mathbf{p}^{\mathbf{b}}$ given by

$$\mathbf{p}_{\mathbf{k}, \mathbf{t}} = \mathbf{p}_{\mathbf{k}}^{\mathbf{b}} = \frac{w_{\mathbf{k}}^{\mathbf{b}}}{\sum_{j=1}^N w_j^{\mathbf{b}}}$$

and draw an action $\mathbf{l}_{\mathbf{t}}$ from $\{1, \dots, N\}$ according to it;

- 2 if $\mathbf{t} = (\mathbf{b} + 1)m$ for some integer \mathbf{b} , perform the update

$$w_{\mathbf{k}}^{\mathbf{b}+1} = w_{\mathbf{k}}^{\mathbf{b}} e^{\eta (\tilde{r}(\mathbf{p}^{\mathbf{b}}, \hat{\Delta}^{\mathbf{b}}))_{\mathbf{k}}} \quad \text{for each } \mathbf{k} = 1, \dots, N,$$

where for all Δ , $\tilde{r}(\cdot, \Delta)$ is a sub-gradient of $\rho(\cdot, \Delta)$.

performance

By optimizing parameters $\mathbf{m} \approx \sqrt{n}$ and $\eta \approx n^{-1/4} \sqrt{\ln N}$, we get the regret bound

$$\begin{aligned} \max_{\mathbf{p}} \rho(\mathbf{p}, \mathbf{H}(\cdot, \bar{\mathbf{q}}_n)) - \frac{1}{n} \sum_{t=1, \dots, n} r(\mathbf{l}_t, \mathbf{y}_t) \\ = O(n^{-1/4} \sqrt{\ln(nN/\delta)}) \end{aligned}$$

which holds with probability $> 1 - \delta$.

the general case

Now the random feedback depends on the action–outcome pairs:
 $\mathbf{H}(\mathbf{I}_t, \mathbf{J}_t)$.

Again, we need to estimate the (unobserved) $\mathbf{H}(\cdot, \bar{\mathbf{q}}_n)$. Let

$$\hat{\mathbf{h}}_{i,t} = \frac{\delta_{s_t}}{\mathbf{p}_{i,t}} \mathbb{I}_{[I_t=i]} .$$

Then $\hat{\mathbf{h}}_{i,t}$ is conditionally unbiased:

$$\mathbb{E}_t \left[\hat{\mathbf{h}}_{i,t} \right] = \frac{\mathbf{1}}{\mathbf{p}_{i,t}} \mathbb{E}_t \left[\delta_{s_t} \mathbb{I}_{[I_t=i]} \right] = \mathbf{H}(i, \mathbf{J}_t) .$$

the general case

Define

$$\hat{\Delta}^b = \Pi \left(\frac{1}{m} \sum_{t=bm+1}^{(b+1)m} [\hat{h}_{i,t}]_{i=1,\dots,N} \right)$$

This will be close to

$$\Delta^b = \frac{1}{m} \sum_{t=bm+1}^{(b+1)m} H(\cdot, J_t)$$

provided $\mathbf{p}_{i,t}$ is not too small.

the strategy

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

- 1 if $bm + 1 \leq t < (b + 1)m$ for some integer b , choose the distribution $\mathbf{p}_t = \mathbf{p}^b = (1 - \gamma)\tilde{\mathbf{p}}^b + \gamma\mathbf{u}$, where $\tilde{\mathbf{p}}^b$ is defined component-wise as

$$\tilde{p}_k^b = \frac{w_k^b}{\sum_{j=1}^N w_j^b}$$

and \mathbf{u} denotes the uniform distribution,

$$\mathbf{u} = (1/N, \dots, 1/N);$$

- 2 draw an action \mathbf{l}_t from $\{1, \dots, N\}$ according to it;
- 3 if $t = (b + 1)m$ for some integer b , perform the update

$$w_k^{b+1} = w_k^b e^{\eta(\tilde{r}(\mathbf{p}^b, \hat{\Delta}^b))_k} \quad \text{for each } k = 1, \dots, N,$$

where for all $\Delta \in \mathcal{F}$, $\tilde{r}(\cdot, \Delta)$ is a sub-gradient of $\rho(\cdot, \Delta)$.

performance

By choosing $\mathbf{m} \approx \mathbf{n}^{3/5}$, $\gamma \approx \mathbf{n}^{-1/5}$, and $\eta \approx \mathbf{n}^{-1/5} \sqrt{\ln \mathbf{N}}$, we obtain

$$\begin{aligned} \max_{\mathbf{p}} \rho(\mathbf{p}, \mathbf{H}(\cdot, \bar{\mathbf{q}}_n)) - \frac{1}{\mathbf{n}} \sum_{t=1, \dots, n} r(\mathbf{l}_t, \mathbf{y}_t) \\ = \mathbf{O}(\mathbf{n}^{-1/5} \mathbf{N} \sqrt{\ln(\mathbf{nN}/\delta)}) \end{aligned}$$

which holds with probability $> 1 - \delta$.

For deterministic feedback this can be improved to $\mathbf{O}(\mathbf{n}^{-1/3} \mathbf{N}^{2/3} \sqrt{\ln(1/\delta)})$.

In the deterministic case the rates (as a function of \mathbf{n}) cannot be improved.

remarks

The strategies involve computation of \mathbf{l}_2 projections to a convex set and computation of (sub)gradients of piecewise linear concave functions.

These can be done in time polynomial in \mathbf{N} and $|\mathbf{S}|$.

Are the obtained rates optimal in the case of random feedback?

hannan consistency

Our strategies are *Hannan consistent* whenever

$$\max_{\mathbf{p}} \rho(\mathbf{p}, \mathbf{H}(\cdot, \mathbf{q})) = \max_{i=1, \dots, N} r(i, \mathbf{q}) .$$

Piccolboni and Schindelhauer (2001), and Cesa-Bianchi, Lugosi, and Stoltz (2006) investigate this special case for deterministic feedback.

For example, Hannan consistency is possible if there is a matrix \mathbf{K} such that $\mathbf{R} = \mathbf{KH}$.

We have a different sufficient condition (independent of the rewards):

If \mathbf{H} doesn't have identical columns, then for all \mathbf{q} ,

$$\max_{\mathbf{p}} \rho(\mathbf{p}, \mathbf{H}(\cdot, \mathbf{q})) = \max_{i=1, \dots, N} r(i, \mathbf{q}) .$$