

Δευτέρα 4/12, Θεωρία Παιγνίων (Βασίλης Τραφόρας)

Ανασκόπηση - Μοντέλα επαναλαμβανόμενης μάθησης (online learning) συνεχούς χρόνου.

- Κάθε χρ. βζιχημή \rightarrow διάνυσμα πληρωμών $u_t \in \mathbb{R}^A$
- \rightarrow Επιλέγει μεκρή στρατηγική $x_t \in \Delta(A)$
- \rightarrow Λαμβάνει $u_t(x_t) = \langle u_t, x_t \rangle$

- Μεταμέτεια (Regret): $Reg(T) = \max_{p \in \Delta(A)} \int_0^T [u_t(p) - u_t(x_t)] dt.$

- Αν ο παίκτης ακολουθεί (EW), τότε: $Reg(T) \leq \log A.$

Online Learning σε διακριτό χρόνο.

Sequence of events Για κάθε χρ. βζιχημή $t=1, 2, \dots$

- Επιλογή μεκρης στρατηγικής $x_t \in \Delta(A).$
- Επιλογή κωδάρης στρατηγικής $a_t \sim x_t.$
- Διάνυσμα πληρωμής: $u_t \in [0, 1]^A.$
- Πληρωμή $u_t(a_t) = v_t, a_t$

Repeat

Μεταμέτεια/Regret: $\sum_{t=1}^T (u_t(a) - u_t(a_t)) = RReg_\alpha(T)$ [Στοχαστικά]. → Realized Regret.

$E[u_t(a)] - E[u_t(a_t)]$ ← Μέση Διαφορά.
 $a_t \sim p \in \Delta(A)$ $a_t \sim x_t \in \Delta(A)$

$$Reg_p(T) = \sum_{t=1}^T \langle u_t, p - x_t \rangle = \sum_{t=1}^T u_t(p) - u_t(x_t)$$

Έχουμε φτάσει στον ορισμό: $Reg_p(T) = \sum_{t=1}^T E[u_t(p) - u_t(a_t)]$

Συνολική μεταμέτεια (Worst-Case Regret):

$Reg(T) = \max_{p \in \Delta(A)} \sum_{t=1}^T \langle u_t, p - x_t \rangle.$

Ποιο το πρόβλημα?

Αν ξεκινήσουμε με τον ορισμό του Realized Regret:

$$RRegret(T) = \sum_{t=1}^T [u_t(a) - u_t(a_t)]$$

Θα έχουμε ως προς μία μεκρή $p \in \Delta(A)$

$$RReg_p(T) = \sum_{t=1}^T [u_t(p) - u_t(a_t)]$$

Οπότε, η χειρότερη περίπτωση θα είναι: $\max_{p \in \Delta(A)} R \text{Reg}(T) = \max_{p \in \Delta(A)} \sum_{t=1}^T [u_t(p) - u_t(a_t)]$

Μέσο Regret = $E_{a_t, \chi_t} \left[\max_{p \in \Delta(A)} \sum_{t=1}^T [u_t(p) - u_t(a_t)] \right]$

• Υπάρχει διαφορά αν κοιτάσουμε:

- ① Των μέσων τιμών ενός μεγίστου (μέσο Regret)
- ② Το μέγιστο μίας μέσων τιμών (pseudo-regret)

Εμείς θα μετρήσουμε αποκλειστικά το 2ο δηλ.:

$$\text{Reg}(T) = \max_{p \in \Delta(A)} \sum_{t=1}^T \langle v_t, p - x_t \rangle$$

ή στην βροχαστική περίπτωση:

$$\overline{\text{Reg}}(T) = \max_{p \in \Delta(A)} E \left[\sum_{t=1}^T \langle v_t, p - x_t \rangle \right]$$

Πληροφόρηση του Παικτη: Αρκετές διαφορετικές περιπτώσεις:

① Πλήρης πληροφόρηση (full information).

Παρατηρούμενη ποσότητα: $v_t \in [0, 1]^T$

② Ευθόρυνες παρατηρήσεις πληροφοριών:

Παρατηρούμενη ποσότητα: $\hat{v}_t = v_t + z_t$

← δείγμα τυχαίου θορύβου.

③ "Bandit" / Παρατήρηση πραγματικής πληροφορίας

Παρατηρούμενη ποσότητα: $u_t(a_t) = v_t, a_t \in [0, 1]$

Παράδειγμα: R-P-S

Παίκτης A (εστιακός) παίζει βροχαστική $\chi_t = (\frac{1}{2}, \frac{1}{3}, \frac{1}{6})$.

Παίκτης B (περιβάλλον): $y_t = (\frac{1}{2}, \frac{1}{2}, 0)$.

$$v_t = \begin{pmatrix} 0 & -1 & 1 \\ 1 & 0 & -1 \\ -1 & 1 & 0 \end{pmatrix} \cdot y_t = \begin{pmatrix} 0 & -1 & 1 \\ 1 & 0 & -1 \\ -1 & 1 & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \\ 0 \end{pmatrix} = \begin{pmatrix} -\frac{1}{2} \\ \frac{1}{2} \\ 0 \end{pmatrix}$$

Περίπτωση 1: R $\begin{pmatrix} -\frac{1}{2} \end{pmatrix}$ P $\begin{pmatrix} \frac{1}{2} \end{pmatrix}$ S $\begin{pmatrix} 0 \end{pmatrix}$.

Περίπτωση 2: $z_R = 0,02$ $z_P = 0,03$ $z_S = 0,5$
 $\begin{pmatrix} -0,10 \end{pmatrix}$ $\begin{pmatrix} 0,23 \end{pmatrix}$ $\begin{pmatrix} -0,5 \end{pmatrix}$

Περίπτωση 3: Ο παίκτης A παίζει

$a_t = (0, 1, 0) \xrightarrow{\chi_t} \text{Paper}$

Άρα παρατηρεί $x = \frac{1}{2} \times$.

Γενική μορφή πληροφορίας πρώτης τάξης των χρονικών βιζυγίων το παίκτη παρατηρεί
 το δείγμα $\hat{v}_t = v_t + u_t + b_t$ όπου: $E[u_t | \mathcal{F}_t] = 0$

τυχαίο βφάλμα
 συστηματικό βφάλμα (bias)

↓ παρεχόν/ιστορία της διαδικασίας μέχρι των βιζυγίων t.

$$\rightarrow b_t = E[\hat{v}_t | \mathcal{F}_t] - v_t$$

Υποθέτουμε: ① Συστηματικό βφάλμα: $\|b_t\| \leq B_t$.

② Διασπορά τυχαίου βφάλματος: $E[\|u_t\|^2 | \mathcal{F}_t] \leq \sigma_t^2$.

③ Δεύτερη ροπή βήματος: $E[\|\hat{v}_t\|^2 | \mathcal{F}_t] \leq \mu_t^2$.

Πλάνο: ① Θα πράξουμε ένα φράγμα μεταμέλειας (Regret Bound) για τον ΕΩ ως προς B_t, σ_t, μ_t .

② Θα το εφαρμόσουμε σε κάθε μοντέλο πληροφορίας.

Αλγόριθμος Hedge / Exponential Weights

$$y_t = v_t$$

Συνεχές χρόνο

$$y_{t+1} = y_t + \gamma v_t$$

Διακριτός χρόνος

$$\chi_t = \Lambda(y_t) = \frac{(\exp(\gamma_1, t), \dots, \exp(\gamma_n, t))}{\sum_{\alpha \in A} \exp(\gamma_\alpha, t)}$$

Τι Regret εφαρμόζει ο Hedge?

Η λογική θα είναι ακριβώς όπως πριν κ' θα βασιστεί στην συνάρτηση δυναμικού.

$$\Phi(y) = \log \sum_{\alpha \in A} \exp(\gamma_\alpha)$$

softmax / log-sum-exp

Βασική Ιδιότητα: $\nabla \Phi(y) = \Lambda(y)$.

Βήμα 1: $y_{t+1} = y_t + \gamma \hat{v}_t$. Θα μελετήσουμε τον αλλαγή $\Phi(y_{t+1})$ σε σύγκριση με $\Phi(y_t)$.

Λήμμα: Έστω $y^+ = y + w$. Τότε, $\Phi(y^+) \leq \Phi(y) + \langle \Lambda(y), w \rangle + \frac{1}{2} \|w\|_\infty^2$ "max over alpha"

Απόδειξη: Από Taylor (με υποβοήθημα Lagrange): $\Phi(y^+) = \Phi(y) + \langle \nabla \Phi(y), y^+ - y \rangle + \frac{1}{2} (y^+ - y)^T \nabla^2 \Phi(y') (y^+ - y)$.

για κάποιο $y' \in [y, y^+]$.

$$\textcircled{I} = \langle \Lambda(y), y^+ - y \rangle = \langle \Lambda(y), w \rangle$$

Για τον \textcircled{II} θα πρέπει να υπολογίσουμε τον Hessian με βεβαιότητα: $\frac{\partial^2 \Phi}{\partial \gamma_\alpha \partial \gamma_\beta} = \frac{\partial}{\partial \gamma_\beta} \frac{\exp(\gamma_\alpha)}{\sum_\gamma \exp(\gamma_\gamma)}$

$$= \dots = \chi_\alpha (\delta_{\alpha\beta} - \chi_\beta)$$

↳ Kronecker $\delta_{\alpha\beta} = \begin{cases} 0, & \alpha \neq \beta \\ 1, & \alpha = \beta \end{cases}$.
α, β ∈ 1, 2, ..., n

$$\textcircled{\text{II}} \Rightarrow \sum_{\alpha, \beta} \chi_\alpha (\delta_{\alpha\beta} - \chi_\beta) \omega_\alpha \omega_\beta$$

$$= \sum_{\alpha} \chi_\alpha \omega_\alpha^2 - \left(\sum_{\alpha} \chi_\alpha \omega_\alpha \right)^2 \leq \sum_{\alpha} \chi_\alpha \omega_\alpha^2 \leq \left(\sum_{\alpha} \chi_\alpha \right) \max_{\beta} \omega_\beta^2 \leq \| \omega \|_\infty^2$$

Βήμα 2: Άρα, $\Phi(y_{t+1}) \leq \Phi(y_t) + \langle \lambda(y_t), y_{t+1} - y_t \rangle + \frac{1}{2} \|y_{t+1} - y_t\|_\infty^2$

$$= \Phi(y_t) + \gamma \langle \chi_t, \hat{v}_t \rangle + \frac{\gamma^2}{2} \|\hat{v}_t\|_\infty^2 = \Phi(y_t) + \gamma \langle \chi_t, \hat{v}_t \rangle + \gamma \langle \chi_t, u_t \rangle + \gamma \langle v_t, p - \chi_t \rangle + \frac{\gamma^2}{2} \|u_t\|_\infty^2$$

$$+ \gamma \langle v_t, p \rangle$$

Cauchy-Schwarz

Βήμα 3: Τυχερόκοποιμε: $\gamma \sum_{t=1}^T \langle v_t, p - \chi_t \rangle \leq \Phi(y_t) - \Phi(y_{t+1}) + 2\gamma \sum_{t=1}^T B_t$

$$\mathbb{E} \left[\sum_{t=1}^T \langle v_t, p - \chi_t \rangle \right] \leq \frac{\Phi(y_t)}{\gamma} + \mathbb{E} \left[\langle y_{t+1}, p \rangle - \Phi(y_{t+1}) \right] + \sum_{t=1}^T B_t$$

$$+ \gamma \sum_{t=1}^T \langle v_t, \chi_t \rangle + \frac{\gamma^2}{2} \sum_{t=1}^T \|\hat{v}_t\|_\infty^2$$

$$+ \gamma \sum_{t=1}^T \langle v_t, p \rangle$$

$$\text{Reg}_p(\tau) + 0 + \frac{\gamma}{2} \sum_{t=1}^T M_t^2$$