

Advanced ML

& AGT

Class 11

Bandits w.

Knapsacks

Bandits with Knapsacks

① מוטיב ציה ומפול

② אלגוריתם מרכזי

③ חוסר תחומן

④ אלגוריתם צמ"ב UCB

מספר מוצרים

d מוצרים
 B_i עותקים של מוצר i

קונה: ערך $U_t = (U_{t,1}, \dots, U_{t,d})$

מוכר: מחירים $P_t = (P_{t,1}, \dots, P_{t,d})$

המוצא:

לקני המוצרים i עבור

$$U_{t,i} \geq P_{t,i}$$

* אפשר גם מוצרים אחרים לקנייה

מוצרים: Dynamic Pricing

רוצים למכור B מוצרים זרים

בזמן t

באזקוח (עם ערך U_t)

המוכר מציג מחיר P_t

אם $U_t \geq P_t$ יש קנייה

עבור $B=D$ (אין מנפחת כמות)

האין בעיה

מה לעשות אם $B < D$?

עלויות להישאר על המוצרים

הבעיה פורמלית

פירמטרים: k פעולות

d משתנים

תקציבים B_1, \dots, B_d

T זמן

בזמן t :

בחרים פעולה $a_t \in [k]$

רווחים תוצאה $O_t = (r_t, c_{t,1}, \dots, c_{t,d}) \in [0,1]^d$

$r_t \in [0,1]$ רווח

$c_{t,i} \in [0,1]$ צריכת משאב i

צריכים שאם ק"מ משאב i עברו

$$\sum c_{t,i} > B_i$$

Bandit with Knapsacks

המטרה: ק"מ אילוצים שלובים
לצביק עק"ם

משאבים בעל כמות מוגבלת
משאב i של כמות B_i

רווחם רווחים עדיפו עמקס

צ'מיקה בפרזמן t

בחרים פעולה a_t

מקבלים רווח r_t

צורכים משאבים

משאב i בכמות $0 \leq c_{t,i}$

המטרה:

מדידת אמינות. את הרמת

$$\sum r_t$$

מסלול נוסף

"37% ב-1 מסלול אחר"

MS

$$\forall i \quad B_i = T$$

MAB '60/60

מבט על אלגוריתם:

אם a הוא הפעולה P_a

$$\theta_t \sim P_a$$

כדי להעריך את הפעולה

$$P \text{ מהפעולה } M_t \in [a, 1]^{k \times (d+1)}$$

הפעולה $M_t(a)$ היא הפעולה a ברגע t

$$M_t(a) = (r_t(a); c_{t,1}(a), \dots, c_{t,d}(a))$$

$$\theta_t = M_t(a_t) \text{ כולל } \theta$$

תמורה צינור אולם צ'ים

T זמן

B תקצ'ים

פאוליס תשלום (מחיר)

מטרה לבצע מקס. משלם

לצורך יש לרכוש מניית v_t

אמס'ים אם $v_t \leq p_t$

משלם קוצ'ים
אחרים

$$\vec{Q}_t = \begin{cases} (1; p, 1) \\ (0; 0, 1) \end{cases}$$

צ'אלמאלי

תמורה צ'ינ

T מספר דקומות

B מספר אורק'ים

פאוליס מחיר'ים

לקונה יש לרכוש v_t

אקונה אם $v_t \geq p_t$

מכר לקונה
לכל לקונה

$$\vec{D}_t = \begin{cases} (p; 1, 0) \\ (0; 0, 1) \end{cases}$$

Repeated auction

T זמן (תקופות)

B איתקם של מוצר

פצולה ממיר מני'מ'ס θ_t

P_t ה'מ'ר'ר בו (מכר ה'מ'וצר

(אם (מכר, ר'מ'ות θ_t)

$$O_t = \begin{cases} (P_t, 1, 0) & \text{מוצר (מכר ב'מ'ר'ר P_t)} \\ (0, 0, 1) & \text{אחרת} \end{cases}$$

Pay per click ad auction

T זמן (מסמ'ס)

B ת'ק'צ'ה של ה'ל'ק'ה

פצולה ד'ה'צ'ה מוצ'ה a

v_a ת'ש'ל'ם = ר'ו'ל

אם ב'ו'צ' ק'י'ק

(ב'מ'ת'ר'ל'ם q_a ל'א'י'צ'ו'ע'ה!)

מ'כ'ר'ה: ד'ס'י'ם א'ר ת'ק'צ'יק ה'ל'ק'ה

$$O_t = \begin{cases} (v_a, v_a, 1) & \text{אם ה'י'ה ק'י'ק} \\ (0, 0, 1) & \text{אחרת} \end{cases}$$

הצגת אל המוצר

(1) צריכת משאבים מקצו (Explore-Exploit)
 עלולה להציק לאחר כך.

(2) הכרוח עלולה - לא המטרה

צריך לחשוב על סכום ההכנסות

(3) הפעולה האובדנית, לא המוצר

ציון: שתי פעולות, תקציב B

$$T = 2B$$

בם זמן t , $d = \{$
 $a_1 \rightarrow (1, 0, 1)$

$a_2 \rightarrow (0, 1, 1)$

כל פעולה בוצרת תשלום B

שימוש בשני הפעולות (החלופין)

יש B!

Repeated Bidding on a Budget

לקנות המלח של קונה

T מספר המכרזים

B תקציב

פעולה b_t (bid המוצר)

$$p_t \leq b_t$$

מטרה: לקנות מספר מקס של מוצרים
 בתקציב B

$$\vec{b}_t = \begin{cases} (1; p_t, 1) & \text{קניית המחר p_t } \\ (0; 0, 1) & \text{לא קניית} \end{cases}$$

פעולות $a \in [k]$

תוצאה $(a=i)$:

$$O_t = (0 \dots 0, \underbrace{CPU_t, mem_t, bw_t}_{i \text{ מסנה}}, 0 \dots 0, 1)$$

\downarrow
מב

מס"מים כאשר אחת המטלות
"מלאה"

Job scheduling

k מספר המכונות

דרכי מכונה i : משאבים:

$CPU_i, memory_i, bandwidth_i$

J מספר עבודות

כל עבודה צורכת משאבים

רואים את הכמות רק אחרי
היזמון של מסוּמ.

מטרה: להביא למקסימום את מספר
העבודות שניתן ליישם.

MAB סטוכסטי:

ALG בומר את הפסוקי
הסוכה ב'ול

MAB w. Knapsacks

נבנה שנספיק לבמן התבטול
על סוכול

Benchmark

$$OPT \triangleq \sup_{ALG} E[REW(ALG|I)]$$

אלגוריתם
פסוקי
הבסיה

BwK

Linear Relaxation: $OPT_{LP} \geq OPT$

Linear Program

Lagrange func.: Game value = OPT_{LP}

Lagrange Game

Learning in Games: Avg Play \approx Nash

Repeated Game

Large reward at stopping time

BwK

Lagrange BwK

אנליזה

Linear Relaxation

לפתור תוכנית ליניארית
התחבולת D פשוט

$$OPT_{LP} \geq \frac{OPT}{T}$$

Lagrange Game

לפתור משחק סכום-אפס בקוץ LP
שחקן אחר בחר פאזה a
שחקן שני בחר משאב b

$$OPT_{LP} = \text{ערך המשחק}$$

BwK

Linear Program

Lagrange Game

Repeated Game

BwK

Linear Relaxation: $DPT_{LP} \geq OPT$

Lagrange func.: Game value = DPT_{LP}

Learning in Games: Avg Play \approx Nash

Large reward at stopping time

Repeated Lagrange Game

משחק חוזר $L_t(a, i)$

התמורה המשחק המקורי $E[L_t] = L$

נפתור ע"י אלג' regret

עברי משחק סכום אדם

(גורם כולשון regret)

Reward at stopping

הכרוך קצת העצירה

צומה ערוך המשחק

(גורם של regret)

הקצמה

$c_i(a)$, $r(a)$ תחומי רוח וזכייה

נתון B_i/B - כח משלוח i ה- $B = \min_{1 \leq i \leq n} B_i$

התקציה של כל משלוח B

הצמ "זרק" הקציה B/T

רוח אכס

mult arm

זרייה אכס, מלכד הצמ

$$OPT_{LP} \geq OPT \quad \text{קיצוץ}$$

הוכחה: ALG הוא אלגוריתם

T של null action פשוט, ALG של μ $\Sigma < T$

ALG של התפלגות אקראית D

$$D(a) = E \left[\frac{1}{T} \sum_{t=1}^T \mathbb{1}\{a_t = a\} \right] = \frac{1}{T} \sum_{t=1}^T P_r[a_t = a]$$

$$E[REW(ALG)] = r(D) \cdot T \quad \text{קיצוץ}$$

$$E[r_t] = \sum_a P_r[a_t = a] E[r_t | a_t = a]$$

$$E[REW] = E \left[\sum_{t=1}^T r_t \right] = \sum_a r(a) \underbrace{\sum_t P_r[a_t = a]}_{T \cdot D(a)} = T \cdot r(D)$$

Linear Relaxation

D התפלגות על Σ

$$r(D) = \sum_a D(a) r(a) \quad C_i(D) = \sum_a D(a) c_i(a)$$

לשאלת D-ה אם כן μ
הוא ו-3 כי זהו הפתרון
הכי טוב

$$\max r(D)$$

$$\text{s.t.} \quad D \in \Delta([k])$$

$$T C_i(D) \leq B \quad \forall i \in [d]$$

$$1 \geq OPT_{LP} \quad \text{הוא LP ה-1}$$

כיוון שהאשורים 318 רק 750 T

$$T \cdot C_i(D) \leq B$$

LP - \Rightarrow D feasible מתוך D פתרון

$$E[\text{REW}(\text{ALG})] = r(D) \leq \text{OPT}_{LP}$$

לכן

*

המשק הכולל:

$$\sum_{t=1}^T E[C_{t,i}] = T \cdot C_i(D) \quad \text{באופן ממוצע}$$

$$E[C_{t,i}] = \sum_a \Pr[a_t = a] E[C_{t,i} | a]$$

$$\sum_{t=1}^T E[C_{t,i}] = \sum_a \underbrace{\sum_{t=1}^T \Pr[a_t = a]}_{T \cdot D(a)} C_i(a)$$

$$= T \cdot C_i(D)$$

דוגמה: עבור שני משתתפים (D^*, λ^*) Nash

$$0 < \lambda_i^* \quad \text{אם } 0 \leq 1 - \frac{I}{B} C(D^*) \quad (B)$$

D^* פתרון ל LP (C)

$$OPT_{LP} = \mathcal{L}(D^*, \lambda^*) \quad (E)$$

הוכחה: עבור Nash משתתפים

$$\mathcal{L}(D^*, \lambda) \geq \mathcal{L}(D^*, \lambda^*) \geq \mathcal{L}(D, \lambda^*)$$

$$\gamma_i = 1 - \frac{I}{B} C_i(D^*) \quad \text{ל } \gamma_i$$

נכנס אל (D) עם משתתפים:

$$\gamma_i \geq 0 \quad \text{עבור } 1 = \lambda_i^* \quad (i)$$

$$(1 > \lambda_i^* \quad \text{כאשר}) \quad \gamma_i \geq 0 \quad (ii)$$

$$\text{במקרה } 0 < \lambda_i^* \quad \gamma_i > 0 \quad (iii)$$

Lagrange פונקציה

LP ה Lagrange פונקציה

$$\mathcal{L}(D, \lambda) = r(D) + \sum_{i \in [d]} \lambda_i \left(1 - \frac{I}{B} C_i(D)\right)$$

לכל משתתף i נכנס:

$$\mathcal{L}(a, i) = r(a) + 1 - \frac{I}{B} C_i(a)$$

משתתף i מקסימום \max פונקציה a

משתתף i מינימום \min פונקציה

$$\min_{\lambda} \max_D \mathcal{L}(D, \lambda) = \max_D \min_{\lambda} \mathcal{L}(D, \lambda) = OPT_{LP}$$

(ii) $\gamma_i \geq 0$ (מספרים חיוביים עבור $\lambda_i < 1$)

כשליטה, נבחר משתנה i עם

$\lambda_i < 1$ - $\gamma_i < 0$ נניח

לפי פילוסוף λ : $\lambda_i = 1$, $\lambda_{i'} = 0$ ($i' \neq i$)

נקבל $\mathcal{L}(D^*, \lambda) < \mathcal{L}(D^*, \lambda^*)$, סתירה!

(iii) $\gamma_i > 0$ כל $\lambda_i = 0$ וכל

נניח $0 < \lambda_i$ - $\gamma_i > 0$

$r(D^*) < \mathcal{L}(D^*, \lambda^*)$

נבחר λ כך : $\lambda(\text{null}) = 1$, $\lambda(a) = 0$ ($a \neq \text{null}$)

$\mathcal{L}(D^*, \lambda) > r(D^*) = \mathcal{L}(D^*, \lambda^*)$

סתירה!

הוכחת הלמה חלק k :

(i) אם $\lambda_i = 1$ אז $\gamma_i = 1 - \frac{T}{B} c_i(D^*) \geq 0$

כל i בהצטרף $\lambda_i = \frac{B}{T} = c_i(D^*)$ - $\gamma_i = 0$

נקבע פשוט a , ולפי ההתאמה D

$$D(\text{null}) = D^*(\text{null}) + D^*(a) \quad D(a') = D^*(a) \quad D(a) = 0$$

$a' \neq a$

$$0 \leq \mathcal{L}(D^*, \lambda^*) - \mathcal{L}(D, \lambda^*)$$

$$= [r(D^*) - r(D)] + \frac{T}{B} [c_i(D^*) - c_i(D)]$$

$$= D^*(a) \left[r(a) - \frac{T}{B} c_i(a) \right]$$

$$\leq D^*(a) \left[1 - \frac{T}{B} c_i(a) \right]$$

נסתכל בהתאמה ונקבל $\gamma_i \geq 0$.

הוכחת תורת

תורת ה

מכאן

LP- \Rightarrow D^* feasible

$$r(D^*) = L(D^*, \lambda^*)$$

כל D feasible

$$r(D^*) = L(D^*, \lambda^*) \geq L(D, \lambda^*) \geq r(D)$$

מכאן D^* optimal LP

$$OPT_{LP} = r(D^*) = L(D^*, \lambda^*)$$

*

Lagrange BwK אלגוריתם

d, k, B, T : קבועים

(bandit בלבד) ALG-max בומר פחולה
(full-info בלבד) ALG-min בומר מטאב

בזמן t :

$a_t \in [k]$ בומר ALG-max : מטקביס

$i_t \in [d]$ בומר ALG-min

פחולה a_t מבוצעת

$\theta_t = (r_t(a_t), c_{t,1}(a_t), \dots)$ קובלים

$\mathcal{L}_t(a_t, i_t)$ קובלה ALG-max

$\forall i \in [d]$ $\mathcal{L}_t(a_t, i)$ קובלה ALG-min

Repeated Lagrange Game

\mathcal{L}_t קובלה M_t קובלה e t בזמן

$$\mathcal{L}_t(a, i) = r_t(a) + 1 - \frac{1}{B} c_{t,i}(a)$$

$$E[\mathcal{L}_t(a, i)] = \mathcal{L}(a, i) \quad \text{קובלה}$$

מטקב

a_t בומר ALG-max פחולה

i_t בומר ALG-min מטאב

$\mathcal{L}_t(a_t, i_t)$ קובלה

קובלה:

יש מטקב קובלה סתמא קובלה
קובלה קובלה

ממשק'ים סוביקטס

$T \in [T]$ משך זמן

(\bar{a}_t, \bar{i}_t) וקטור הממוצע

של δ_T -Nash

$$\sum \delta_T = R_1(T) + R_2(T) + \text{err}_T$$

$$\text{err}_T = \left| \sum_{t \in [T]} L_t(i_t, j_t) - L(i_t, j_t) \right|$$

כיוון

$$\forall i \in [d] \quad L(\bar{a}_T, i) \geq \text{OPT}_{LP} - \delta_T$$

הוכחה:

$R_1(T)$ is regret of ALG-max

$$\bar{a}_T = \frac{1}{T} \sum_{t \in [T]} a_t \rightarrow \text{average}$$

$R_2(T)$ is regret of ALG-min

$$\bar{i}_T = \frac{1}{T} \sum_{t \in [T]} i_t \rightarrow \text{average}$$

: Hoefding 3/17/2018

$$\text{err}_{\tau,i} \leq \sqrt{TK \log \frac{dT}{\delta}} \quad \forall \tau \in [T], i \in [d]$$

מה שגורם לregret יחסית קטן

$$R_1(\tau) \leq z \sqrt{TK \log(T/\delta)}$$

$$R_2(\tau) \leq z \sqrt{T \log(dT/\delta)}$$

רצוננו

$$\underbrace{\text{OPT-REW}}_{\text{regret}} \leq O\left(\frac{T}{B} \cdot \sqrt{TK \log(dT/\delta)}\right)$$

$$B = \Omega(T) \text{ קטן יותר פחות}$$

$$B = \Omega(\sqrt{T}) \text{ קטן יותר פחות}$$

רצוננו

$$\mathcal{L}(\bar{a}_{\tau,i}) \geq \text{OPT}_{LP} - \delta_{\tau} \quad \forall i$$

$$\tau \delta_{\tau} = R_1(\tau) + R_2(\tau) + \text{err}_{\tau}$$

$$\text{REW} \geq \tau \mathcal{L}(\bar{a}_{\tau,i}) + (T-\tau) \text{OPT}_{LP} - \text{err}_{\tau}$$

$$\text{REW} \geq T \cdot \text{OPT}_{LP} - R_1(\tau) - R_2(\tau) - \text{err}_{\tau}$$

: (ההנחה) פחות

$$\text{err}_{\tau} \leq z \sqrt{TK \log \left(\frac{dT}{\delta}\right)}$$

$$z = \max |\mathcal{L}(a_i)| \leq \frac{T}{B}$$

$$\mathcal{L}(a_i) \in \left[1 - \frac{T}{B}, 2\right]$$

חסם פתרון $\min\{OPT, OPT\sqrt{K/B}\}$

הרעיון: צומה ל MAB

הרווח תמיד 1

הצרכים $\{0, B\}$

$Pr[c(a^*)=1] = \frac{1-\epsilon}{2}$ פאזה אופט

$Pr[c(a)=1] = \frac{1}{2}$ פאזה אחרת

הזמן $4B = T$ (לא מיליוני!)

תוחלת OPT $2B / (1-\epsilon)$

הצהרה: היה בצד N MAB

הזמן הוא N

הרווחים דטרמיניסטיים

חסם פתרון (בקיבוץ)

של N:

$$\text{Regret} = \Omega\left(\min\left\{OPT, \sqrt{K}OPT + OPT\sqrt{\frac{K}{B}}\right\}\right)$$

חסם פתרון $\min\{OPT, \sqrt{K}OPT\}$

ל MAB - N

ONLINE אלגוריתם T אלגוריתם
 משתקם פחות אוכל ϵ

אלגוריתם X_1 ϵ -OPT

אלגוריתם X_2 ONLINE ϵ

אלגוריתם ϵ regret

$$\frac{X_1}{(1-\epsilon)^{1/2}} - \frac{X_2}{(1-\epsilon)^{1/2}} = \Theta(X_1 - X_2) = \Theta(\epsilon T)$$

$$\text{regret} = \Theta\left(\sqrt{\frac{k}{B}} \text{OPT}\right) = \Theta(\epsilon T) \quad \text{אם}$$

עמוד מספר MAB

$$T \leq c \frac{k}{\epsilon^2}, \quad T = 2B - \sqrt{B} \log \frac{1}{\epsilon}$$

הסתברות שיהיה לא נצטרך לפני T

צביר פרופורציה $1/3$

כמה זמן t

בתחילת הזמן (שחקן פחותה לא אולי)

ואם מהמספר התחיל MAB!

כמה T ההפרש בצרכי

של OPT ONLINE

$$\Theta(\epsilon \cdot T)$$

$$\epsilon = \Theta\left(\sqrt{\frac{k}{B}}\right) \quad \text{אם } \Theta\left(\sqrt{\frac{k}{B}} \text{OPT}\right) \text{ קבוע}$$

$$a_t = \arg \max_a \text{UCB}_t(a)$$

$$\text{UCB}_t(a_t) \geq \text{UCB}_t(a^*) \geq \mu^*$$

$$\text{Regret} = \sum_{t=1}^T \mu^* - r_t$$

$$\leq \sum_{t=1}^T \text{UCB}_t(a_t) - r_t$$

$$\leq \sqrt{KT}$$

(23'p2) UCB - 1N3 7e5/c

UCB 88 601, 62, N

$$\forall i, t \quad \text{UCB}_t(a_i) \geq \mu_i \geq \text{LCB}_t(a_i)$$

$$\left| \sum_{t=1}^T \text{UCB}_t(a_t) - r_t \right|$$

$$\leq \sum_a \left| \sum_{t: a_t=a} \text{UCB}_t(a_t) - \text{LCB}_t(a) \right|$$

$$\leq \sqrt{KT}$$

$$\delta = \log \frac{kTd}{\epsilon} \quad \epsilon = \sqrt{\frac{\delta k}{B}} + \frac{\delta k_i}{B} \log T \text{ : צד שמאל}$$

$$\text{Regret} = \tilde{O}\left(\text{OPT} \sqrt{\frac{k}{B}} + \sqrt{k \text{OPT}} + k\right)$$

מלבד שם הסתמך הסתמך!

הצבה: הסתמך הצבה

$$UCB_t(r) \geq r \quad LCB_t(a) \leq c$$

$$\left| \sum_{t=1}^T UCB_t(r) \cdot D_t - r_t \right| = O(\sqrt{\delta k \sum r_t})$$

$$\left| \sum_{t=1}^T LCB_t(a) D_t - c_t \right| \leq \epsilon \cdot B \vec{1}$$

(... סתמך ...)

LP(r, a, ε) - הצבה

$$\max_D r(D)$$

$$\text{s.t. } QD \leq B(1-\epsilon) \vec{1}$$

$$D \in \Delta([k])$$

$$C = \begin{bmatrix} -c_1 \\ \vdots \\ -c_d \end{bmatrix}$$

: D ∈ הצבה

$$D_t \leftarrow \text{LP}(UCB_t(r), LCB_t(a), \epsilon)$$

$$\forall i \ a(i) \geq b(i) \iff \vec{a} \succeq \vec{b} \text{ : o.m.f.}$$

$$\begin{aligned}
\text{REW(Alg)} &= \sum_{t=1}^T r_t \\
&= \left[\sum_t \text{UCB}_t(r) D_t \right] - \left[\sum_t \text{UCB}_t(r) D_t - r_t \right] \\
&\geq (1-\epsilon) \text{OPT}_{LP} - \sqrt{\gamma K \cdot \text{REW}} - \delta m
\end{aligned}$$

$$\text{Regret} = O\left(\underbrace{\epsilon \cdot \text{OPT}}_{\frac{\sqrt{\gamma K} \epsilon}{B} \text{OPT}} + \sqrt{\gamma K \text{OPT}} + \gamma K\right)$$

$$\frac{1}{T} \sum_{t=1}^T \text{UCB}_t(r) \cdot D_t = \frac{1}{T} \sum_{t=1}^T \text{LP}(\text{UCB}_t(r), \text{LCB}_t(c), \epsilon)$$

$$\geq \text{LP}(r, c, \epsilon)$$

$$\geq (1-\epsilon) \text{OPT}_{LP}$$

$$\frac{1}{T} \sum_{t=1}^T \text{LCB}_t(c) D_t \leq (1-\epsilon) B \vec{1}$$

$$\frac{1}{T} \sum_{t=1}^T c_t \leq B \vec{1} \quad \text{for }$$

T מוסר גבול לא מוגבל

*

מקורות

גורם ההרצאה

Slivkins, 10 בינואר

חוסם תחתון

Badanidiyuru, Kleinberg, Slivkins, 2018
(Bandits with knapsacks)

אלגוריתם UCB

Agrawal, Devanur, 2014

סיכום ההרצאה

מנסה

Lagrange

מתחם

הקצאה:

חוסם תחתון

אלגוריתם UCB