

Μαθηματικές Ιδιότητες της Στοχαστικής Προσέγγισης

Καθολική προσέγγιση, στοχαστική προσέγγιση και multi-armed bandits

Νίκος Σταμάτης

22 Ιανουαρίου 2021

1. Καθολική Προσέγγιση
2. Στοχαστική Προσέγγιση
3. Multi-armed Bandits

Καθολική Προσέγγιση

Ορισμός

Μια συνάρτηση $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ καλείται σιγμοειδής, εαν

$$\sigma(t) \longrightarrow \begin{cases} 0, & \text{για } t \rightarrow -\infty, \\ 1, & \text{για } t \rightarrow +\infty. \end{cases}$$

Ορισμός

Μια συνάρτηση $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ καλείται *σιγμοειδής*, εαν

$$\sigma(t) \longrightarrow \begin{cases} 0, & \text{για } t \rightarrow -\infty, \\ 1, & \text{για } t \rightarrow +\infty. \end{cases}$$

Παραδείγματα

- Σιγμοειδής: $\sigma(x) = \frac{1}{1+e^{-x}}$.
- Step function: $s(x) = \begin{cases} 0, & x < 0, \\ 1, & x \geq 0. \end{cases}$

Ορισμός

Μια συνάρτηση $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ καλείται σιγμοειδής, εαν

$$\sigma(t) \longrightarrow \begin{cases} 0, & \text{για } t \rightarrow -\infty, \\ 1, & \text{για } t \rightarrow +\infty. \end{cases}$$

Ορισμός

Μια συνάρτηση $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ καλείται *σιγμοειδής*, εαν

$$\sigma(t) \longrightarrow \begin{cases} 0, & \text{για } t \rightarrow -\infty, \\ 1, & \text{για } t \rightarrow +\infty. \end{cases}$$

Ορισμός

Μια συνάρτηση $A : X \rightarrow Y$ μεταξύ δύο διανυσματικών χώρων ονομάζεται *αφφινική*, εαν ισχύει ότι $A(\sum_{i=1}^n \lambda_i x_i) = \sum_{i=1}^n \lambda_i A(x_i)$ για κάθε $n \in \mathbb{N}$, $x_i \in X$ και κάθε $\lambda_i \in \mathbb{R}$ με $\sum_{i=1}^n \lambda_i = 1$.

Ορισμός

Μια συνάρτηση $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ καλείται *σιγμοειδής*, εαν

$$\sigma(t) \longrightarrow \begin{cases} 0, & \text{για } t \rightarrow -\infty, \\ 1, & \text{για } t \rightarrow +\infty. \end{cases}$$

Ορισμός

Μια συνάρτηση $A : X \rightarrow Y$ μεταξύ δύο διανυσματικών χώρων ονομάζεται *αφφινική*, εαν ισχύει ότι $A(\sum_{i=1}^n \lambda_i x_i) = \sum_{i=1}^n \lambda_i A(x_i)$ για κάθε $n \in \mathbb{N}$, $x_i \in X$ και κάθε $\lambda_i \in \mathbb{R}$ με $\sum_{i=1}^n \lambda_i = 1$.

Νευρωνικό δίκτυο, ονομάζουμε κάθε συνάρτηση s της μορφής $s(x) = Tx$, όπου

$$T = A_{m+1} S_m A_m \cdots A_2 S_1 A_1$$

είναι ένας τελεστής ο οποίος ορίζεται απο τις διαδοχικές συνθέσεις αφφινικών μετασχηματισμών $A_i : \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ με σιγμοειδείς συναρτήσεις S_i . Ο αριθμός m μετρά το πλήθος των στρωμάτων του δικτύου, ενώ οι αριθμοί d_i το πλήθος των κόμβων που εμφανίζονται σε κάθε στρώμα.

Ερώτημα

Για ποιες συναρτήσεις ενεργοποίησης σ ισχύει ότι το σύνολο

$$\Sigma_n(\sigma) = \left\{ f: I_n \rightarrow \mathbb{R} : f(x) = \sum_{j=1}^N a_j \sigma(w_j^T x + \theta_j) : \right. \\ \left. N \in \mathbb{N}, a_j \in \mathbb{R}, w_j \in \mathbb{R}^n, \theta_j \in \mathbb{R} \right\} \\ = \text{span} \{ f: f(x) = \sigma(w^T x + \theta) \text{ για } w \in \mathbb{R}^n, \theta \in \mathbb{R} \}$$

είναι πυκνό στον $C(I_n)$;

Θεώρημα (Cybenko, 1989)

Για κάθε συνεχή σιγμοειδή συνάρτηση ενεργοποίησης σ , το σύνολο $\Sigma_n(\sigma)$ είναι πυκνό στον $C(I_n)$.

$$\begin{aligned}\Sigma_n(\sigma) = & \left\{ f: I_n \rightarrow \mathbb{R} : f(x) = \sum_{j=1}^N a_j \sigma(w_j^T x + \theta_j) : \right. \\ & \left. N \in \mathbb{N}, a_j \in \mathbb{R}, w_j \in \mathbb{R}^n, \theta_j \in \mathbb{R} \right\} \\ = & \text{span} \{ f: f(x) = \sigma(w^T x + \theta) \text{ για } w \in \mathbb{R}^n, \theta \in \mathbb{R} \} .\end{aligned}$$

Θεώρημα (Cybenko, 1989)

Για κάθε συνεχή σιγμοειδή συνάρτηση ενεργοποίησης σ , το σύνολο $\Sigma_n(\sigma)$ είναι πυκνό στον $C(I_n)$.

Θεώρημα (Chen, Chen, Liu, 1991)

Κατασκευαστική απόδειξη του Θεωρήματος Cybenko.

$$\begin{aligned}\Sigma_n(\sigma) = & \left\{ f: I_n \rightarrow \mathbb{R} : f(x) = \sum_{j=1}^N a_j \sigma(w_j^T x + \theta_j) : \right. \\ & \left. N \in \mathbb{N}, a_j \in \mathbb{R}, w_j \in \mathbb{R}^n, \theta_j \in \mathbb{R} \right\} \\ = & \text{span} \{ f: f(x) = \sigma(w^T x + \theta) \text{ για } w \in \mathbb{R}^n, \theta \in \mathbb{R} \} .\end{aligned}$$

Θεώρημα (Cybenko, 1989)

Για κάθε συνεχή σιγμοειδή συνάρτηση ενεργοποίησης σ , το σύνολο $\Sigma_n(\sigma)$ είναι πυκνό στον $C(I_n)$.

Θεώρημα (Chen, Chen, Liu, 1991)

Κατασκευαστική απόδειξη του Θεωρήματος Cybenko.

Θεώρημα (Leshno, Lin, Pinkus, Schocken, 1993)

Το σύνολο $\Sigma_n(\sigma)$ είναι πυκνό στον $C(\mathbb{R}^n)$, εαν και μόνο εαν η συνάρτηση ενεργοποίησης σ δεν είναι πολυώνυμο.

$$\begin{aligned}\Sigma_n(\sigma) = & \left\{ f: I_n \rightarrow \mathbb{R} : f(x) = \sum_{j=1}^N a_j \sigma(w_j^T x + \theta_j) : \right. \\ & \left. N \in \mathbb{N}, a_j \in \mathbb{R}, w_j \in \mathbb{R}^n, \theta_j \in \mathbb{R} \right\} \\ = & \text{span} \{ f: f(x) = \sigma(w^T x + \theta) \text{ για } w \in \mathbb{R}^n, \theta \in \mathbb{R} \} .\end{aligned}$$

Θεώρημα (Maïoron, Pinkus, 1999)

Υπάρχει λεία, σιγμοειδής συνάρτηση ενεργοποίησης σ , τέτοια ώστε για κάθε $d \in \mathbb{N}$, κάθε συμπαγές $K \subseteq \mathbb{R}^d$, κάθε $f \in C(K)$ και $\varepsilon > 0$, να υπάρχουν πραγματικές σταθερές $d_i, c_{ij}, \theta_{ij}, \gamma_i$ και διανύσματα $w_{ij} \in \mathbb{R}^d$, με

$$\left| f(x) - \sum_{i=1}^{6d+3} d_i \sigma \left(\sum_{j=1}^{3d} c_{ij} \sigma(w_{ij}^T x + \theta_{ij}) + \gamma_i \right) \right| < \varepsilon$$

για κάθε $x \in K$.

Θεώρημα (Maïorov, Pinkus, 1999)

Υπάρχει λεία, σιγμοειδής συνάρτηση ενεργοποίησης σ , τέτοια ώστε για κάθε $d \in \mathbb{N}$, κάθε συμπαγές $K \subseteq \mathbb{R}^d$, κάθε $f \in C(K)$ και $\varepsilon > 0$, να υπάρχουν πραγματικές σταθερές $d_i, c_{ij}, \theta_{ij}, \gamma_i$ και διανύσματα $w_{ij} \in \mathbb{R}^d$, με

$$\left| f(x) - \sum_{i=1}^{6d+3} d_i \sigma \left(\sum_{j=1}^{3d} c_{ij} \sigma(w_{ij}^T x + \theta_{ij}) + \gamma_i \right) \right| < \varepsilon$$

για κάθε $x \in K$.

Θεώρημα (Kolmogorov-Arnold, 1957)

Υπάρχουν σταθερές $\lambda_1, \dots, \lambda_d$ τέτοιες ώστε $\sum_{j=1}^d \lambda_j \leq 1$, και συνεχείς συναρτήσεις $\phi_1, \dots, \phi_{2d+1}$ από το $[0, 1]$ στον εαυτό του, με την ιδιότητα ότι κάθε $f \in C[0, 1]^d$ μπορεί να γραφτεί ως

$$f(x_1, \dots, x_d) = \sum_{i=1}^{2d+1} g \left(\sum_{j=1}^d \lambda_j \phi_i(x_j) \right),$$

όπου $g \in C[0, 1]$ μια συνάρτηση που εξαρτάται από την f .

Στοχαστική Προσέγγιση

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.
- Να βρεθεί θ με $M(\theta) = a$.

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.
- Να βρεθεί θ με $M(\theta) = a$.

Εργαλεία:

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.
- Να βρεθεί θ με $M(\theta) = a$.

Εργαλεία:

- Μπορούμε να προσομοιώνουμε από την $H(y | x)$ για κάθε τιμή του x .

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.
- Να βρεθεί θ με $M(\theta) = a$.

Εργαλεία:

- Μπορούμε να προσομοιώνουμε από την $H(y | x)$ για κάθε τιμή του x .

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.
- Να βρεθεί θ με $M(\theta) = a$.

Εργαλεία:

- Μπορούμε να προσομοιώνουμε από την $H(y | x)$ για κάθε τιμή του x .

Λύση:

Πλαίσιο:

- Μεταβλητή απόκρισης Y , επεξηγηματική μεταβλητή X .
- Δοθέντος ότι $X = x$, η κατανομή της $Y(x)$ είναι $P[Y(x) \leq y] = H(y | x)$.
- Έστω $M(x)$ η μέση τιμή της Y δοθέντος του x , δηλ. $M(x) = \mathbb{E}[Y | X = x]$.
πχ. στο γραμμικό μοντέλο $Y = Xb + \varepsilon$, έχουμε $E[Y | X = x] = x'b$.
- Να βρεθεί θ με $M(\theta) = a$.

Εργασία:

- Μπορούμε να προσομοιώνουμε από την $H(y | x)$ για κάθε τιμή του x .

Λύση:

- Μέσω της αναδρομικής ακολουθίας

$$x_{n+1} = x_n + a_n(a - y_n),$$

όπου η παρατήρηση y_n έχει προσομοιωθεί από την $H(y | x_n)$.

Θεώρημα (Robbins-Monro 1951)

Αν υπάρχει $C > 0$ τ.ω. $x \in \mathbb{R}$, $P[|Y(x)| \leq C] = 1$, αν η $M(x) = \mathbb{E}[Y | X = x]$ είναι αύξουσα με $M(\theta) = a$ και $M'(\theta) > 0$, τότε για κάθε $(a_n)_n \in \ell_2^+ \setminus \ell_1$ η $(x_n)_n$ που κατασκευάζεται από τον αλγόριθμο Robbins-Monro συγκλίνει κατά πιθανότητα στο θ .

Θεώρημα (Robbins-Monro 1951)

Αν υπάρχει $C > 0$ τ.ω. $x \in \mathbb{R}$, $P[|Y(x)| \leq C] = 1$, αν η $M(x) = \mathbb{E}[Y | X = x]$ είναι αύξουσα με $M(\theta) = a$ και $M'(\theta) > 0$, τότε για κάθε $(a_n)_n \in \ell_2^+ \setminus \ell_1$ η $(x_n)_n$ που κατασκευάζεται από τον αλγόριθμο Robbins-Monro συγκλίνει κατά πιθανότητα στο θ .

ΑΛΓΟΡΙΘΜΟΣ ROBBINS-MONRO

ΒΗΜΑ 1 Επιλέγουμε $(a_n)_n \in \ell_2^+ \setminus \ell_1$ και $x_1 \in \mathbb{R}$.

ΒΗΜΑ 2 Υποθέτουμε ότι έχουμε κατασκευάσει το x_1, \dots, x_{n-1} . Προσομοιώνουμε μια παρατήρηση y_n από την κατανομή $H(y | x_n)$ και θέτουμε $x_{n+1} = x_n + a_n(a - y_n)$.

ΒΗΜΑ 3 Επιστρέφουμε στο Βήμα 2.

Ορισμός

Έστω X σύνολο και $H : X \rightarrow X$ συνάρτηση. Ένα σημείο $x_0 \in X$ καλείται σταθερό σημείο της H , εαν $H(x_0) = x_0$.

Ορισμός

Έστω X σύνολο και $H : X \rightarrow X$ συνάρτηση. Ένα σημείο $x_0 \in X$ καλείται σταθερό σημείο της H , εαν $H(x_0) = x_0$.

Ορισμός

Μια συνάρτηση $H : (X, \rho) \rightarrow (X, \rho)$ ορισμένη σε έναν μετρικό χώρο (X, ρ) , καλείται β -συστολή, όπου $\beta > 0$ σταθερά, εαν

$$\rho(Hx, Hy) \leq \beta \rho(x, y) \text{ για κάθε } x, y \in X.$$

Ορισμός

Έστω X σύνολο και $H : X \rightarrow X$ συνάρτηση. Ένα σημείο $x_0 \in X$ καλείται σταθερό σημείο της H , εαν $H(x_0) = x_0$.

Ορισμός

Μια συνάρτηση $H : (X, \rho) \rightarrow (X, \rho)$ ορισμένη σε έναν μετρικό χώρο (X, ρ) , καλείται β -συστολή, όπου $\beta > 0$ σταθερά, εαν

$$\rho(Hx, Hy) \leq \beta \rho(x, y) \text{ για κάθε } x, y \in X.$$

Θεώρημα σταθερού σημείου του Banach (1922)

Έστω (X, ρ) πλήρης μετρικός χώρος και $H : X \rightarrow X$ συνάρτηση η οποία είναι β -συστολή για κάποιο $\beta \in (0, 1)$. Τότε η H έχει μοναδικό σταθερό σημείο x_0 . Επιπλέον, η ακολουθία $(x_n)_n$ που ορίζεται αναδρομικά ως $x_1 = y$ και $x_{n+1} = H(x_n)$ για κάθε $n \in \mathbb{N}$, συγκλίνει στο x_0 για οποιαδήποτε αρχική επιλογή $y \in X$.

Ερώτημα

Έστω X σύνολο και $H : X \rightarrow X$ συνάρτηση με μοναδικό σταθερό σημείο το $x_0 \in X$. Πώς μπορούμε να προσδιορίσουμε το σταθερό της σημείο, όταν δε γνωρίζουμε τις ακριβείς τιμές της H ;

Ερώτημα

Έστω X σύνολο και $H : X \rightarrow X$ συνάρτηση με μοναδικό σταθερό σημείο το $x_0 \in X$. Πώς μπορούμε να προσδιορίσουμε το σταθερό της σημείο, όταν δε γνωρίζουμε τις ακριβείς τιμές της H ;

Αλγόριθμος

Κατασκευάζουμε αναδρομικά την ακολουθία:

$$x_{n+1} = (1 - \gamma_n)x_n + \gamma_n(Hx_n + w_n),$$

όπου $(\gamma_n)_n$ είναι κατάλληλη ακολουθία στο $(0, 1]$, και οι όροι w_n αντιστοιχούν στο “θόρυβο” των τιμών Hx_n .

Ψευδοσυστολές

Έστω $(X, \|\cdot\|)$ χώρος με νόρμα. Μια συνάρτηση $H : X \rightarrow X$ καλείται *ψευδοσυστολή*, εαν υπάρχουν $x^* \in X$ και $\beta \in [0, 1)$ τέτοια ώστε

$$\|Hx - Hx^*\| \leq \beta \|x - x^*\|$$

για κάθε $x \in X$.

Ψευδοσυστολές

Έστω $(X, \|\cdot\|)$ χώρος με νόρμα. Μια συνάρτηση $H : X \rightarrow X$ καλείται *ψευδοσυστολή*, εαν υπάρχουν $x^* \in X$ και $\beta \in [0, 1)$ τέτοια ώστε

$$\|Hx - Hx^*\| \leq \beta \|x - x^*\|$$

για κάθε $x \in X$.

Μονότονοι τελεστές

Έστω $(X, \|\cdot\|, \leq)$ χώρος με νόρμα εφοδιασμένος με μια γραμμική διάταξη. Μια συνάρτηση $H : X \rightarrow X$ καλείται *μονότονη*, εαν $Hx \leq Hy$ για κάθε $x, y \in X$ με $x \leq y$.

Πρόταση

Έστω $(r_n)_n$ η ακολουθία που ορίζεται από την αναδρομική σχέση

$$r_{n+1} = (1 - \gamma_n)r_n + \gamma_n(H_n r_n + w_n + u_n),$$

όπου

1. η ακολουθία $(\gamma_n)_n$ είναι τέτοια ώστε $\sum_{n=1}^{\infty} \gamma_n(i) = \infty$ και $\sum_{n=1}^{\infty} \gamma_n(i)^2 < \infty$ για κάθε $i = 1, \dots, N$.
2. Η ακολουθία $(w_n)_n$ έχει την ιδιότητα ότι

$$\mathbb{E}[w_n(i) \mid \mathcal{F}_n] = 0 \quad \text{και} \quad \mathbb{E}[w_n(i)^2 \mid \mathcal{F}_n] \leq A + B\|r_n\|^2.$$

3. Κάθε H_n είναι ψευδοσυστολή ως προς την ίδια νόρμα $\|\cdot\|_{\xi}$, με το ίδιο σταθερό σημείο r^* και την ίδια σταθερά $\beta \in [0, 1)$.
4. Υπάρχει ακολουθία μη-αρνητικών τυχαίων μεταβλητών $(\theta_n)_n$ η οποία συγκλίνει στο μηδέν σχεδόν παντού, τέτοια ώστε

$$\|u_n\|_{\infty} \leq \theta_n (1 + \|r_n\|_{\xi})$$

για κάθε $n \in \mathbb{N}$.

Τότε η $(r_n)_n$ συγκλίνει στο r^* σχεδόν παντού.

Πρόταση

Έστω $(r_n)_n$ η ακολουθία που ορίζεται από την αναδρομική σχέση

$$r_{n+1} = (1 - \gamma_n)r_n + \gamma_n(Hr_n + w_n),$$

όπου

1. η ακολουθία $(\gamma_n)_n$ είναι τέτοια ώστε $\sum_{n=1}^{\infty} \gamma_n(i) = \infty$ και $\sum_{n=1}^{\infty} \gamma_n(i)^2 < \infty$ για κάθε $i = 1, \dots, N$.
2. Η ακολουθία $(w_n)_n$ έχει την ιδιότητα ότι

$$\mathbb{E}[w_n(i) \mid \mathcal{F}_n] = 0 \quad \text{και} \quad \mathbb{E}[w_n(i)^2 \mid \mathcal{F}_n] \leq A + B\|r_n\|^2.$$

3. Για τον τελεστή H ισχύει ότι

3.1 είναι μονότονος, δηλαδή $Hx \leq Hy$ για κάθε $x \leq y$.

3.2 Για κάθε $\lambda > 0$ και $r \in \mathbb{R}^N$, ισχύει ότι:

$$Hr - \lambda e \leq H(r - \lambda e) \leq H(r + \lambda e) \leq Hr + \lambda e, \text{ όπου } e = (1, \dots, 1).$$

3.3 Έχει μοναδικό σταθερό σημείο, $Hr^* = r^*$.

Εαν η $(r_n)_n$ είναι φραγμένη σχεδόν παντού, τότε συγκλίνει στο r^* σχεδόν παντού.

- MDP με σύνολο καταστάσεων S .

- MDP με σύνολο καταστάσεων S .
- Σε κάθε στάδιο n του προβλήματος, παρατηρούμε την τρέχουσα κατάσταση $i \in S$ της διαδικασίας και επιλέγουμε μια απόφαση $a \in A(i)$.

- MDP με σύνολο καταστάσεων S .
- Σε κάθε στάδιο n του προβλήματος, παρατηρούμε την τρέχουσα κατάσταση $i \in S$ της διαδικασίας και επιλέγουμε μια απόφαση $a \in A(i)$.
- Η διαδικασία μεταβαίνει στην κατάσταση j με πιθανότητα $p_{ij}(a)$. Οι πιθανότητες μετάβασης εξαρτώνται από την τρέχουσα κατάσταση i και την απόφαση a που λήφθηκε.

- MDP με σύνολο καταστάσεων S .
- Σε κάθε στάδιο n του προβλήματος, παρατηρούμε την τρέχουσα κατάσταση $i \in S$ της διαδικασίας και επιλέγουμε μια απόφαση $a \in A(i)$.
- Η διαδικασία μεταβαίνει στην κατάσταση j με πιθανότητα $p_{ij}(a)$. Οι πιθανότητες μετάβασης εξαρτώνται από την τρέχουσα κατάσταση i και την απόφαση a που λήφθηκε.
- Κάθε μετάβαση επιφέρει κόστος $c(i, a, j)$.

- MDP με σύνολο καταστάσεων S .
- Σε κάθε στάδιο n του προβλήματος, παρατηρούμε την τρέχουσα κατάσταση $i \in S$ της διαδικασίας και επιλέγουμε μια απόφαση $a \in A(i)$.
- Η διαδικασία μεταβαίνει στην κατάσταση j με πιθανότητα $p_{ij}(a)$. Οι πιθανότητες μετάβασης εξαρτώνται από την τρέχουσα κατάσταση i και την απόφαση a που λήφθηκε.
- Κάθε μετάβαση επιφέρει κόστος $c(i, a, j)$.
- Αξία πολιτικής π : Το αναμενόμενο κόστος $J^\pi(i)$ όταν η διαδικασία ξεκινάει από την κατάσταση $i_0 = i$, και ακολουθείται η πολιτική $\pi = (\mu_0, \mu_1, \dots)$,

$$J^\pi(i) = \lim_{N \rightarrow \infty} \mathbb{E} \left[\sum_{m=0}^N c(i_m, \mu_m(i_m), i_{m+1}) \mid i_0 = i \right].$$

- MDP με σύνολο καταστάσεων S .
- Σε κάθε στάδιο n του προβλήματος, παρατηρούμε την τρέχουσα κατάσταση $i \in S$ της διαδικασίας και επιλέγουμε μια απόφαση $a \in A(i)$.
- Η διαδικασία μεταβαίνει στην κατάσταση j με πιθανότητα $p_{ij}(a)$. Οι πιθανότητες μετάβασης εξαρτώνται από την τρέχουσα κατάσταση i και την απόφαση a που λήφθηκε.
- Κάθε μετάβαση επιφέρει κόστος $c(i, a, j)$.
- Αξία πολιτικής π : Το αναμενόμενο κόστος $J^\pi(i)$ όταν η διαδικασία ξεκινάει από την κατάσταση $i_0 = i$, και ακολουθείται η πολιτική $\pi = (\mu_0, \mu_1, \dots)$,

$$J^\pi(i) = \lim_{N \rightarrow \infty} \mathbb{E} \left[\sum_{m=0}^N c(i_m, \mu_m(i_m), i_{m+1}) \mid i_0 = i \right].$$

- Συνάρτηση βέλτιστης τιμής $J^*(i)$: Το ελάχιστο δυνατό κόστος κάτω από οποιαδήποτε πολιτική, όταν η διαδικασία ξεκινάει από την κατάσταση i ,

$$J^*(i) = \min_{\pi} J^\pi(i).$$

Εξίσωση Bellman

$$J^*(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + J^*(j)) \right\}.$$

Εξίσωση Bellman

$$J^*(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + J^*(j)) \right\}.$$

Μέθοδος διαδοχικών προσεγγίσεων

Ξεκινάμε από κάποια αυθαίρετη συνάρτηση J_0 , και εν συνεχεία ορίζουμε αναδρομικά

$$J_{n+1}(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + J_n(j)) \right\}.$$

Εξίσωση Bellman

$$J^*(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + J^*(j)) \right\}.$$

Μέθοδος διαδοχικών προσεγγίσεων

Ξεκινάμε από κάποια αυθαίρετη συνάρτηση J_0 , και εν συνεχεία ορίζουμε αναδρομικά

$$J_{n+1}(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + J_n(j)) \right\}.$$

Η σύγκλιση της ακολουθίας $(J_n)_n$ στην J^* εξασφαλίζεται μέσω του θεωρήματος σταθερού σημείου του Banach, καθώς ο τελεστής $T : C(S) \rightarrow C(S)$, που ορίζεται ως

$$(Tf)(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + f(j)) \right\}, \quad f \in C(S), \quad i \in S,$$

αποτελεί συστολή.

Ερώτημα

Μπορούμε να λύσουμε την εξίσωση Bellman, όταν δε γνωρίζουμε τις πιθανότητες μετάβασης $p_{ij}(a)$ και τα κόστη $c(i, a, j)$;

Ερώτημα

Μπορούμε να λύσουμε την εξίσωση Bellman, όταν δε γνωρίζουμε τις πιθανότητες μετάβασης $p_{ij}(a)$ και τα κόστη $c(i, a, j)$;

Q-Παράγοντες

Για $(i, a) \in S \times A(i)$, ορίζουμε τον βέλτιστο Q-παράγοντα $Q^*(i, a)$ ως

$$Q^*(0, a) = 0 \text{ and } Q^*(i, a) = \sum_{j=0}^N p_{ij}(a) (c(i, a, j) + J^*(j)) \text{ για } i = 1, \dots, N,$$

όπου J^* η συνάρτηση βέλτιστης τιμής.

Ερώτημα

Μπορούμε να λύσουμε την εξίσωση Bellman, όταν δε γνωρίζουμε τις πιθανότητες μετάβασης $p_{ij}(a)$ και τα κόστη $c(i, a, j)$;

Q-Παράγοντες

Για $(i, a) \in S \times A(i)$, ορίζουμε τον βέλτιστο Q-παράγοντα $Q^*(i, a)$ ως

$$Q^*(0, a) = 0 \text{ and } Q^*(i, a) = \sum_{j=0}^N p_{ij}(a) (c(i, a, j) + J^*(j)) \text{ για } i = 1, \dots, N,$$

όπου J^* η συνάρτηση βέλτιστης τιμής.

Ο Q^* ικανοποιεί την εξίσωση

$$Q^*(i, a) = \sum_{j=0}^N p_{ij}(a) \left(c(i, a, j) + \min_{b \in A(j)} Q^*(j, b) \right) \text{ for } i = 1, \dots, N,$$

- Εξίσωση Bellman για την J^* :

$$J^*(i) = \min_{a \in A(i)} \left\{ \sum_{j \in S} p_{ij}(a) (c(i, a, j) + J^*(j)) \right\}.$$

- Εξίσωση Bellman για τους Q-παράγοντες:

$$Q^*(i, a) = \sum_{j=0}^N p_{ij}(a) \left(c(i, a, j) + \min_{b \in A(j)} Q^*(j, b) \right).$$

Ο ΑΛΓΟΡΙΘΜΟΣ Q-LEARNING

- ΒΗΜΑ 1** Επιλέγουμε συναρτήσεις $(\gamma_n)_n$ ορισμένες στο σύνολο $\tilde{S} = \{(i, a) : i = 1, \dots, N, a \in A(i)\}$, τ.ω. $\sum_{n=0}^{\infty} \gamma_n(i, a) = \infty$ και $\sum_{n=0}^{\infty} \gamma_n(i, a)^2 < \infty$ για κάθε $i = 1, \dots, N$ και $a \in A(i)$.
- ΒΗΜΑ 2** Αρχικοποιούμε με τη συνάρτηση $Q_0(i, a)$ για $i = 1, \dots, N$ και $a \in A(i)$.
- ΒΗΜΑ 3** Υποθέτουμε πως έχουμε δημιουργήσει τις τιμές $Q_n(i, a)$ για κάθε ζεύγος (i, a) για κάποιο $n \geq 0$. Για κάθε ζεύγος (i, a) , προσομοιώνουμε μια παρατήρηση $j_{i,a}$ από την κατανομή $p_{i,\cdot}(a)$ και θέτουμε $Q_{n+1}(i, a) = (1 - \gamma_n(i, a))Q_n(i, a) + \gamma_n(i, a) \left(c(i, a, j_{i,a}) + \min_{b \in A(j_{i,a})} Q_n(j_{i,a}, b) \right)$.
- ΒΗΜΑ 4** Επιστρέφουμε στο ΒΗΜΑ 3.

Θεώρημα (Watkins 1989, Tsitsiklis 1994)

Θεωρούμε την ακολουθία $(Q_n)_n$ που ορίζεται αναδρομικά ως

$$Q_{n+1}(i, a) = (1 - \gamma_n(i, a)) Q_n(i, a) + \gamma_n(i, a) \left(c(i, a, j) + \min_{b \in A(j)} Q_n(j, b) \right),$$

όπου η κατάσταση j έχει προσομοιωθεί από την κατανομή $p_{i,\cdot}(a)$ και $(\gamma_n)_n$ είναι ακολουθία με την ιδιότητα ότι $\sum_{n=0}^{\infty} \gamma_n(i, a) = \infty$ και $\sum_{n=0}^{\infty} \gamma_n(i, a)^2 < \infty$ για κάθε $i = 1, \dots, N$ και $a \in A(i)$. Αν όλες οι πολιτικές είναι γνήσιες, τότε $Q_n(i, a) \rightarrow Q^*(i, a)$ για κάθε $i, a \in A(i)$ σχεδόν παντού.

Θεώρημα (Watkins 1989, Tsitsiklis 1994)

Θεωρούμε την ακολουθία $(Q_n)_n$ που ορίζεται αναδρομικά ως

$$Q_{n+1}(i, a) = (1 - \gamma_n(i, a)) Q_n(i, a) + \gamma_n(i, a) \left(c(i, a, j) + \min_{b \in A(j)} Q_n(j, b) \right),$$

όπου η κατάσταση j έχει προσομοιωθεί από την κατανομή $p_{i,\cdot}(a)$ και $(\gamma_n)_n$ είναι ακολουθία με την ιδιότητα ότι $\sum_{n=0}^{\infty} \gamma_n(i, a) = \infty$ και $\sum_{n=0}^{\infty} \gamma_n(i, a)^2 < \infty$ για κάθε $i = 1, \dots, N$ και $a \in A(i)$. Αν όλες οι πολιτικές είναι γνήσιες, τότε $Q_n(i, a) \rightarrow Q^*(i, a)$ για κάθε $i, a \in A(i)$ σχεδόν παντού.

Ο τελεστής $H : C(\tilde{S}) \rightarrow C(\tilde{S})$ που ορίζεται ως

$$(HQ)(i, a) = \sum_{j=0}^N p_{ij}(a) \left(c(i, a, j) + \min_{b \in A(j)} Q(j, b) \right).$$

για Q στο $C(\tilde{S})$, προκύπτει είτε ότι είναι συστολή ως προς κάποια κατάλληλη νόρμα, είτε ότι είναι μονότονος.

Multi-armed Bandits

- Δύο διαφορετικοί πληθυσμοί με κατανομές κερδών p_A, p_B και αναμενόμενες τιμές $\mu_A < \mu_B$ αντίστοιχα.

- Δύο διαφορετικοί πληθυσμοί με κατανομές κερδών p_A, p_B και αναμενόμενες τιμές $\mu_A < \mu_B$ αντίστοιχα.
- Σε κάθε γύρο επιλέγουμε έναν πληθυσμό και τραβάμε μια τιμή από αυτόν, εισπράττοντάς την ως κέρδος.

- Δύο διαφορετικοί πληθυσμοί με κατανομές κερδών p_A, p_B και αναμενόμενες τιμές $\mu_A < \mu_B$ αντίστοιχα.
- Σε κάθε γύρο επιλέγουμε έναν πληθυσμό και τραβάμε μια τιμή από αυτόν, εισπράττοντάς την ως κέρδος.
- *Ερώτημα:* Μπορούμε να βρούμε μια στρατηγική που μακροπρόθεσμα θα μας δίνει αναμενόμενο κέρδος όσο πιο κοντά γίνεται στο μ_B ;

Θεώρημα Robbins (1952)

Ασυμπτωτικά γίνεται να επιτύγχουμε αναμενόμενο κέρδος ακριβώς ίσο με μ_B , σχεδόν βεβαίως.

Ο Κανόνας του Robbins

Θεώρημα Robbins (1952)

Ασυμπτωτικά γίνεται να επιτύχουμε αναμενόμενο κέρδος ακριβώς ίσο με μ_B , σχεδόν βεβαίως.

Ο ΚΑΝΟΝΑΣ ΤΟΥ ROBBINS

- ΒΗΜΑ 1** Επιλέγουμε $J_A, J_B \subseteq \mathbb{N}$ άπειρα με μηδενική πυκνότητα.
- ΒΗΜΑ 2** Υποθέτουμε ότι έχουμε τραβήξει δείγμα x_1, \dots, x_{n-1} στους πρώτους $n - 1$ γύρους. Αν $n \in J_A$, τραβάμε την X_n από τον πληθυσμό A . Αν $n \in J_B$, τραβάμε την X_n από τον B .
- ΒΗΜΑ 3** Αν $n \notin J_A \cup J_B$, θέτουμε $a_n = \frac{\sum_{\{i: X_i \sim F_A\}} x_i}{\#\{i: X_i \sim F_A\}}$ και $b_n = \frac{\sum_{\{i: X_i \sim F_B\}} x_i}{\#\{i: X_i \sim F_B\}}$. Αν $a_n \geq b_n$, τραβάμε την X_n από τον πληθυσμό A , διαφορετικά από τον B .

Ο Κανόνας του Robbins

Θεώρημα Robbins (1952)

Ασυμπτωτικά γίνεται να επιτύγχουμε αναμενόμενο κέρδος ακριβώς ίσο με μ_B , σχεδόν βεβαίως.

Ο ΚΑΝΟΝΑΣ ΤΟΥ ROBBINS

- ΒΗΜΑ 1** Επιλέγουμε $J_A, J_B \subseteq \mathbb{N}$ άπειρα με μηδενική πυκνότητα.
- ΒΗΜΑ 2** Υποθέτουμε ότι έχουμε τραβήξει δείγμα x_1, \dots, x_{n-1} στους πρώτους $n - 1$ γύρους. Αν $n \in J_A$, τραβάμε την X_n από τον πληθυσμό A . Αν $n \in J_B$, τραβάμε την X_n από τον B .
- ΒΗΜΑ 3** Αν $n \notin J_A \cup J_B$, θέτουμε $a_n = \frac{\sum_{\{i: X_i \sim F_A\}} x_i}{\#\{i: X_i \sim F_A\}}$ και $b_n = \frac{\sum_{\{i: X_i \sim F_B\}} x_i}{\#\{i: X_i \sim F_B\}}$. Αν $a_n \geq b_n$, τραβάμε την X_n από τον πληθυσμό A , διαφορετικά από τον B .

Μηδενική πυκνότητα:

$$d(J_A) = \lim_{n \rightarrow \infty} \frac{\#J_A \cap \{1, \dots, n\}}{n} = 0.$$

- Έχουμε K το πλήθος bandits με κατανομές $f(x, \theta_1), \dots, f(x, \theta_K)$ και μέσες τιμές μ_1, \dots, μ_K . Έστω $\mu^* = \max\{\mu_i\}_i$.

- Έχουμε K το πλήθος bandits με κατανομές $f(x, \theta_1), \dots, f(x, \theta_K)$ και μέσες τιμές μ_1, \dots, μ_K . Έστω $\mu^* = \max\{\mu_i\}_i$.
- Συνολικό κέρδος S_n μετά απο n γύρους.

- Έχουμε K το πλήθος bandits με κατανομές $f(x, \theta_1), \dots, f(x, \theta_K)$ και μέσες τιμές μ_1, \dots, μ_K . Έστω $\mu^* = \max\{\mu_i\}_i$.
- Συνολικό κέρδος S_n μετά απο n γύρους.
- Regret $R_n(\theta) = n\mu^* - \mathbb{E}[S_n]$.

- Έχουμε K το πλήθος bandits με κατανομές $f(x, \theta_1), \dots, f(x, \theta_K)$ και μέσες τιμές μ_1, \dots, μ_K . Έστω $\mu^* = \max\{\mu_i\}_i$.
- Συνολικό κέρδος S_n μετά απο n γύρους.
- Regret $R_n(\theta) = n\mu^* - \mathbb{E}[S_n]$.
- Robbins: $\frac{R_n(\theta)}{n} = \mu^* - \frac{\mathbb{E}[S_n]}{n} \rightarrow 0$.

- Έχουμε K το πλήθος bandits με κατανομές $f(x, \theta_1), \dots, f(x, \theta_K)$ και μέσες τιμές μ_1, \dots, μ_K . Έστω $\mu^* = \max\{\mu_i\}_i$.
- Συνολικό κέρδος S_n μετά απο n γύρους.
- Regret $R_n(\theta) = n\mu^* - \mathbb{E}[S_n]$.
- Robbins: $\frac{R_n(\theta)}{n} = \mu^* - \frac{\mathbb{E}[S_n]}{n} \rightarrow 0$.

Το Θεώρημα του Robbins εγγυάται ότι ναι μεν $R_n(\theta) \rightarrow \infty$, αλλά η σύγκλιση είναι πιο αργή από την ακολουθία n , δηλ. $R_n(\theta) = o(n)$.

- Έχουμε K το πλήθος bandits με κατανομές $f(x, \theta_1), \dots, f(x, \theta_K)$ και μέσες τιμές μ_1, \dots, μ_K . Έστω $\mu^* = \max\{\mu_i\}_i$.
- Συνολικό κέρδος S_n μετά απο n γύρους.
- Regret $R_n(\theta) = n\mu^* - \mathbb{E}[S_n]$.
- Robbins: $\frac{R_n(\theta)}{n} = \mu^* - \frac{\mathbb{E}[S_n]}{n} \rightarrow 0$.

Το Θεώρημα του Robbins εγγυάται ότι ναι μεν $R_n(\theta) \rightarrow \infty$, αλλά η σύγκλιση είναι πιο αργή από την ακολουθία n , δηλ. $R_n(\theta) = o(n)$.

Μπορούμε να επιτύχουμε ακόμα πιο αργή σύγκλιση για το regret;

Ορισμός (Kullback-Leibler divergence)

Για τις κατανομές $f(x; \lambda)$ and $f(x; \mu)$ ορίζουμε την απόσταση Kullback - Leibler $I(\lambda, \mu)$ ως

$$I(\lambda, \mu) = \int f(x; \lambda) \ln \frac{f(x; \lambda)}{f(x; \mu)} dx.$$

Ορισμός (Kullback-Leibler divergence)

Για τις κατανομές $f(x; \lambda)$ and $f(x; \mu)$ ορίζουμε την απόσταση Kullback - Leibler $I(\lambda, \mu)$ ως

$$I(\lambda, \mu) = \int f(x; \lambda) \ln \frac{f(x; \lambda)}{f(x; \mu)} dx.$$

Υποθέσεις (Υ1-Υ3):

Ορισμός (Kullback-Leibler divergence)

Για τις κατανομές $f(x; \lambda)$ and $f(x; \mu)$ ορίζουμε την απόσταση Kullback - Leibler $I(\lambda, \mu)$ ως

$$I(\lambda, \mu) = \int f(x; \lambda) \ln \frac{f(x; \lambda)}{f(x; \mu)} dx.$$

Υποθέσεις (Υ1-Υ3):

- **Συνθήκη Κατανομών:** Μονοπαραμετρική οικογένεια κατανομών για τα rewards, $(f(x; \theta))_{\theta \in \Theta}$ με $\Theta \subseteq \mathbb{R}$.

Ορισμός (Kullback-Leibler divergence)

Για τις κατανομές $f(x; \lambda)$ and $f(x; \mu)$ ορίζουμε την απόσταση Kullback - Leibler $I(\lambda, \mu)$ ως

$$I(\lambda, \mu) = \int f(x; \lambda) \ln \frac{f(x; \lambda)}{f(x; \mu)} dx.$$

Υποθέσεις (Υ1-Υ3):

- **Συνθήκη Κατανομών:** Μονοπαραμετρική οικογένεια κατανομών για τα rewards, $(f(x; \theta))_{\theta \in \Theta}$ με $\Theta \subseteq \mathbb{R}$.
- **Συνθήκη Συνέχειας:** Για κάθε $\lambda, \theta \in \Theta$, αν $(\lambda_n)_n$ τ.ω. $\mu(\lambda_n) \downarrow \mu(\lambda)$, τότε $I(\theta, \lambda_n) \rightarrow I(\theta, \lambda)$.

Ορισμός (Kullback-Leibler divergence)

Για τις κατανομές $f(x; \lambda)$ and $f(x; \mu)$ ορίζουμε την απόσταση Kullback - Leibler $I(\lambda, \mu)$ ως

$$I(\lambda, \mu) = \int f(x; \lambda) \ln \frac{f(x; \lambda)}{f(x; \mu)} dx.$$

Υποθέσεις (Υ1-Υ3):

- **Συνθήκη Κατανομών:** Μονοπαραμετρική οικογένεια κατανομών για τα rewards, $(f(x; \theta))_{\theta \in \Theta}$ με $\Theta \subseteq \mathbb{R}$.
- **Συνθήκη Συνέχειας:** Για κάθε $\lambda, \theta \in \Theta$, αν $(\lambda_n)_n$ τ.ω. $\mu(\lambda_n) \downarrow \mu(\lambda)$, τότε $I(\theta, \lambda_n) \rightarrow I(\theta, \lambda)$.
- **Συνθήκη Πυκνότητας:** Για κάθε $\lambda \in \Theta$, υπάρχει $(\lambda_n)_n$ in Θ τ.ω. $(\mu(\lambda_n))_n$ γνησίως φθίνουσα με $\mu(\lambda_n) \downarrow \mu(\lambda)$.

Θεώρημα Lai-Robbins (1985)

Οι κατανομές των χεριών ικανοποιούν τις υποθέσεις Υ1-Υ3. Για έναν αλγόριθμο που ικανοποιεί την ιδιότητα ότι

$$R_n(\theta) = o(n^a)$$

για κάθε $\theta = (\theta_1, \dots, \theta_K) \in \Theta^K$ και $a > 0$, ισχύει ότι

$$\liminf_n \frac{R_n(\theta)}{\ln n} \geq \sum_{i: \mu(\theta_i) < \mu^*} \frac{\mu^* - \mu(\theta_i)}{I(\theta_i, \theta^*)} > 0,$$

για τα θ για τα οποία τα $\mu(\theta_i)$ δεν είναι όλα ίσα.

Υποθέτουμε ότι υπάρχουν

- συναρτήσεις $\tilde{\mu}_n(j)$ που παίζουν το ρόλο του δειγματικού μέσου του j -πληθυσμού, μετά από n -γύρους.

Υποθέτουμε ότι υπάρχουν

- συναρτήσεις $\tilde{\mu}_n(j)$ που παίζουν το ρόλο του δειγματικού μέσου του j -πληθυσμού, μετά από n -γύρους.
- Συναρτήσεις $U_n(j)$ που παίζουν το ρόλο του άνω άκρου ενός διαστήματος εμπιστοσύνης για το μέσο του j -πληθυσμού, μετά από n -γύρους.

Ο Αλγόριθμος Lai-Robbins

Υποθέτουμε ότι υπάρχουν

- συναρτήσεις $\tilde{\mu}_n(j)$ που παίζουν το ρόλο του δειγματικού μέσου του j -πληθυσμού, μετά από n -γύρους.
- Συναρτήσεις $U_n(j)$ που παίζουν το ρόλο του άνω άκρου ενός διαστήματος εμπιστοσύνης για το μέσο του j -πληθυσμού, μετά από n -γύρους.

Ο ΑΛΓΟΡΙΘΜΟΣ LAI-ROBBINS

ΒΗΜΑ 1 Για $m = 1, \dots, K$ επιλέγουμε $\phi(m) = m$.

ΒΗΜΑ 2 Κατά τον $n + 1 \equiv j \bmod K$ γύρο, θέτουμε $I_n = \{m \in \{1, \dots, K\} : T_n(m) \geq \delta n\}$, $j_n = \arg \max \{\tilde{\mu}_n(m) : m \in I_n\}$ και $\tilde{\mu}_n(j_n) = \max \{\tilde{\mu}_n(m) : m \in I_n\}$.

ΒΗΜΑ 3 Αν $\tilde{\mu}_n(j_n) \leq U_n(j)$, τότε $\phi(n + 1) = j$, διαφορετικά $\phi(n + 1) = j_n$.

ΒΗΜΑ 4 Επιστρέφουμε στο ΒΗΜΑ 2.

- Οι συναρτήσεις $\tilde{\mu}_n(j)$ είναι ακριβώς οι δειγματικοί μέσοι του j -πληθυσμού μετά από n -γύρους.

Ο Αλγόριθμος Auer-Bianchi-Fischer

- Οι συναρτήσεις $\tilde{\mu}_n(j)$ είναι ακριβώς οι δειγματικοί μέσοι του j -πληθυσμού μετά από n -γύρους.
- Οι συναρτήσεις $U_n(j)$ δίνονται από την έκφραση

$$\bar{x}_{j,n_j} + \sqrt{\frac{3 \ln n}{2n_j}},$$

όπου \bar{x}_{j,n_j} είναι ο δειγματικός μέσος του j -πληθυσμού και n_j το πλήθος των φορών που επιλέχθηκε ο j -πληθυσμός κατά τους πρώτους n γύρους.

Ο ΑΛΓΟΡΙΘΜΟΣ UPPER CONFIDENCE BOUND

ΒΗΜΑ 1 Επιλέγουμε κάθε πληθυσμό μία φορά.

ΒΗΜΑ 2 Στον $n + 1$ -γύρο, επιλέγουμε τον πληθυσμό που μεγιστοποιεί την έκφραση $\bar{x}_{j,n_j} + \sqrt{\frac{3 \ln n}{2n_j}}$.

ΑΛΓΟΡΙΘΜΟΣ LAI-ROBBINS

ΑΛΓΟΡΙΘΜΟΣ UCB

ΑΛΓΟΡΙΘΜΟΣ LAI-ROBBINS

Συγκεκριμένες υποθέσεις για τις κατανομές των κερδών (παραμετρικές, απαιτήσεις συνέχειας ως προς τη μετρική Kullback-Leibler, κλπ.).

ΑΛΓΟΡΙΘΜΟΣ UCB

Μόναδική υπόθεση είναι οι κατανομές να φέρονται στο $[0, 1]$.

ΑΛΓΟΡΙΘΜΟΣ LAI-ROBBINS	ΑΛΓΟΡΙΘΜΟΣ UCB
Συγκεκριμένες υποθέσεις για τις κατανομές των κερδών (παραμετρικές, απαιτήσεις συνέχειας ως προς τη μετρική Kullback-Leibler, κλπ.).	Μόναδική υπόθεση είναι οι κατανομές να φέρονται στο $[0, 1]$.
Η κατασκευή των δ.ε. είναι δύσκολη υπόθεση. Ακόμα και όταν τα δ.ε. δίνονται, ο υπολογισμός τους είναι απαιτητικός.	Η κατασκευή των δ.ε. είναι εύκολη και υπολογιστικά αποδοτική.

ΑΛΓΟΡΙΘΜΟΣ LAI-ROBBINS	ΑΛΓΟΡΙΘΜΟΣ UCB
Συγκεκριμένες υποθέσεις για τις κατανομές των κερδών (παραμετρικές, απαιτήσεις συνέχειας ως προς τη μετρική Kullback-Leibler, κλπ.).	Μόναδική υπόθεση είναι οι κατανομές να φέρονται στο $[0, 1]$.
Η κατασκευή των δ.ε. είναι δύσκολη υπόθεση. Ακόμα και όταν τα δ.ε. δίνονται, ο υπολογισμός τους είναι απαιτητικός.	Η κατασκευή των δ.ε. είναι εύκολη και υπολογιστικά αποδοτική.
Λογαριθμική απώλεια ασυμπτωτικά.	Λογαριθμική απώλεια ομοιόμορφα.

ΑΛΓΟΡΙΘΜΟΣ LAI-ROBBINS	ΑΛΓΟΡΙΘΜΟΣ UCB
Συγκεκριμένες υποθέσεις για τις κατανομές των κερδών (παραμετρικές, απαιτήσεις συνέχειας ως προς τη μετρική Kullback-Leibler, κλπ.).	Μόναδική υπόθεση είναι οι κατανομές να φέρονται στο $[0, 1]$.
Η κατασκευή των δ.ε. είναι δύσκολη υπόθεση. Ακόμα και όταν τα δ.ε. δίνονται, ο υπολογισμός τους είναι απαιτητικός.	Η κατασκευή των δ.ε. είναι εύκολη και υπολογιστικά αποδοτική.
Λογαριθμική απώλεια ασυμπτωτικά.	Λογαριθμική απώλεια ομοιόμορφα.
Η λογαριθμική σταθερά ισούται με $\frac{1}{2\Delta_j}$ για τους μη βέλτιστους πληθυσμούς j .	Η αντίστοιχη λογαριθμική σταθερά ισούται με $\frac{6}{\Delta_j} > \frac{1}{2\Delta_j}$.



D. Bertsekas and J. Tsitsiklis.
Neuro-Dynamic Programming.
Athena Scientific, 1996.



G. Cybenko.
Approximation by superposition of a sigmoidal function.
Mathematics of Control, Signal, and Systems, 2:303–314, 1989.



T. Lai and H. Robbins.
Asymptotically efficient adaptive allocation rules.
Advances in Applied Mathematics, 2:4–22, 1985.



T. Lattimore and C. Szepesvári.
Bandit Algorithms.
Cambridge University Press, 2020.



P. F. P. Auer, N. Cesa-Bianchi.
Finite-time analysis of the multiarmed bandit problem.
Machine Learning, 47:235–256, 2002.