

کلاس ... منوی 49411448

تمرین کلاسی 2:

نمای این بحث تئوری ROBBINS-MONRO یا اتونیم ROBBINS-MONRO است.
 یک تئوری بدون یکیم و سو آن را از یک یکیم.
 فرض کنید که یک FMPP داریم با فضای state برابر با \mathcal{X} و فضای A و تمیزات
 حالت انتقال P و پاداز r .
 نحوه آتیمیت Q به صورت زیر است:

$$Q_{t+1} = Q_t + \alpha_t [r_t - Q_t]$$

که برای Q_{long} به صورت خاص داریم (حالت خاص از رابطه بالا می شود)

$$Q_{t+1}(x_t, a_t) = Q_t(x_t, a_t) + \alpha_t(x_t, a_t) [I]$$

$$I = r_t + V_{max}(Q_t(x_{t+1}, b)) - Q_t(x_t, a_t) \quad b \in A$$

برای شرایطی داریم که لازم است برای (x, a) عنصر $\mathcal{X} \times A$ (فرض داریم این مجموعه)

$$(1) \sum_t \alpha_t(x, a) < \infty$$

رابطه متناهی برقرار باشد:

$$(2) \sum_t \alpha_t^2(x, a) < \infty$$

$$0 \leq \alpha_t(x, a) \leq 1$$

تئوری جدید نیز این است که یک ترانه فضای Δ_t به صورت خاص است

$$\Delta_{t+1}(u) = (1 - \alpha_t(u)) \Delta_t(u) + \alpha_t(u) F_t(u)$$

و به صورت دیگری می توان نوشت:

$$1) \sum_t \alpha_t^2(u) < \infty, \sum_t \alpha_t(u) < \infty, 0 \leq \alpha_t \leq 1$$

$$2) \|E[F_t(u) | F_t]\|_w \leq \gamma \| \Delta_t \|_w, \text{ with } \gamma < 1$$

$$3) \text{var}[F_{t+1} | F_t] \leq C(1 + \| \Delta_t \|_w^2) \text{ for } C \geq 0$$

Date:

Sub:

۹۹۶۱۱۶۷۸

سید سجاد، فردی

۲- تشریح:

$$Q_{t+1}(x_t, a_t) = (1 - \alpha_t(x_t, a_t)) Q_t(x_t, a_t) + \alpha_t(x_t, a_t) [r_t + \gamma \max_{b \in A} Q_t(x_{t+1}, b)]$$

سازگار، آلمی، تکرار

$$\Delta_t(x_t, a_t) = Q_t(x_t, a_t) - Q^*(x_t, a_t)$$

داریم:

$$\Delta_t(x_t, a_t) = (1 - \alpha_t(x_t, a_t)) \Delta_t(x_t, a_t) + \alpha_t(x_t, a_t) [II]$$

$$II = r_t + \gamma \max_{b \in A} Q_t(x_{t+1}, b) - Q^*(x_t, a_t)$$

تکرار، تکرار، تکرار

$$F_t(x_t, a_t) = r(x_t, a_t, X(x_t, a_t)) + \gamma \max_{b \in A} Q_t(y, b) - Q^*(x_t, a_t)$$

در حالتی که $X(x_t, a_t)$ یک گزینه ممکن است از $X(x_t, a_t)$ و غیره باشد.

$$E[F_t(x_t, a_t) | F_t] = \sum_{j \in \mathcal{A}} p_a(x_t, y) [r(x_t, a_t, y) + \gamma \max_{b \in A} Q_t(y, b) - Q^*(x_t, a_t)]$$

$$= (HQ_t)(x_t, a_t) - Q^*(x_t, a_t)$$

با این روش

$$E[F_t(x_t, a_t) | F_t] = (HQ_t)(x_t, a_t) - HQ^*(x_t, a_t)$$

Date:

Sub:

4951154
سید محمد نوری

$$(Hq)_{(x,a)} = \sum_{y \in X} p_a(x,y) [r(x,a,y) + \gamma \sum_{b \in A} Q_t(y,b)]$$

$$\|Hq_1 - Hq_2\|_{\infty} \leq \gamma \|q_1 - q_2\|_{\infty}$$

بنابراین:

$$\|E[F_t(x,a) | F_t]\|_{\infty} \leq \gamma \|Q_t - Q^*\|_{\infty} = \gamma \|\Delta_t\|_{\infty}$$

و در نتیجه داریم:

$$\text{var}[F_t(x,a) | F_t]$$

$$= E \left[\left(r(x,a, X(x,a)) + \gamma \sum_{b \in A} Q_t(y,b) - Q^*(x,a) - (HQ_t)_{(x,a)} + Q^*(x,a) \right)^2 \right]$$

$$= E \left[\left(r(x,a, X(x,a)) + \gamma \sum_{b \in A} Q_t(y,b) - (HQ_t)_{(x,a)} \right)^2 \right]$$

$$= \text{var} \left[r(x,a, X(x,a)) + \gamma \sum_{b \in A} Q_t(y,b) | F_t \right]$$

$$C = \gamma$$

و در نتیجه داریم:

$$\text{var}[F_t(x,a) | F_t] \leq C (1 + \|\Delta_t\|_{\infty}^2)$$

بنابراین متغیری Δ_t به صورت متراکم $Q^* \approx Q_t$ متراکم می شود.

و متغیری Δ_t متراکم می شود.