

Subject :

Date

بہم خدا

ترتیب ہفتم سری دوم

سوال ۱۱ [۱] Q-value . دیدیم کہ برابر انڈیسی بیٹ ۱ Q-value و V-value
 را پیدا کنیم . اگر فوالم بیٹ صید را پیدا کنیم ، اگر V-value استفادہ کنیم داریم :

$$q(s) = \argmax_a T(s, a, s') [R(s, a, s') + \gamma \bar{v}(s')]$$

کہ در بیٹات پیدا آہن P و R دانستہ باشیم و لہذا چون بی فوالم model-free

این مقادیر را اندازیم . برابر صواب آئینہ زبہ Sample ہی زیادہ داریم کہ ہن model based (محدود)

و لہذا اگر از Q-value استفادہ کنیم ، بہا گریں توانیم بیٹ را بہ روشن زیر اوریاں کنیم :

$$q(s) = \argmax_a Q(s, a)$$

کہ برابر صواب Q (ہا) ، Sample گریں داریم کہ دیگر از مقدار T (P و N) بہینہ زیر استفادہ و
 model free نہ باشیم ، (صواب صمد Q ہا با Q-learning صورت در پذیرا)

است که بهای صورت محل را کند که در ع در صد مواقع (ع ب رکود) یک حرکت

مرکزی برائے (Q. 150) ہائیکہ تا (آن مدت آوردہ انجم دہد)

سبب در اصل رابطه $exPloitation$ و $exPloitation$ با درهم آمیخته و

خوب اگر گفت به روشی Policy iteration of f Policy دارد این است

سرمیت بالاتر داد و از ۳۰ به ۲۵ ریال یاد کرد که استقامت را بد

۱۰ حال اگر اندک در استقامت کنیم به جرات آن زندگانی روش و حال را پس

Policy iteration اتفاق کر اگر state جدید رفاهت شود در این سلسله

مقادیر α و β چون دیگر به explorations بستگی دارند در همان سبک

ما نفعی که ممکن است در حالت جدید آید،

● Greedy: این خریّت را دارد که چهار تا (حق) را (هر چند کم) به ما اجازه

• now in Exploration

(ج) طلب فرض شد داریم:

$$\forall s: E_{\pi'}[Q^{\pi'}(s, a)] \geq E_{\pi}[Q^{\pi}(s, a)]$$

$$\begin{aligned} \forall s: & \sum_{s'} P(s' | s, \pi'(s)) [R(s, \pi'(s), s') + \gamma \max_{a'} Q^{\pi'}(s', a')] \\ & \geq \sum_{s'} P(s' | s, \pi(s)) [R(s, \pi(s), s') + \gamma \max_{a'} Q^{\pi}(s', a')] \quad (*) \end{aligned}$$

حال در هر دو معادله داریم:

$$v^{\pi}(s') = \max_{a'} Q^{\pi}(s', a')$$

$$\begin{aligned} (*) \Rightarrow \forall s: & \sum_{s'} P(s' | s, \pi'(s)) [R(s, \pi'(s), s') + \gamma v^{\pi'}(s')] \\ & \geq \sum_{s'} P(s' | s, \pi(s)) [R(s, \pi(s), s') + \gamma v^{\pi}(s')] \quad (***) \end{aligned}$$

همین طلب فرمول v-value برای سیاست ثابت (Policy Fixed) داریم:

$$v^{\pi} = \sum_{s'} P(s' | s, \pi(s)) [R(s, \pi(s), s') + \gamma v^{\pi}(s')]$$

$$(***) \Rightarrow \forall s: v^{\pi'}(s) \geq v^{\pi}(s)$$