

فرادرس

فراتر از یک کلاس درس
www.faradars.org

آموزش یادگیری تقویتی عمیق در پایتون پیاده‌سازی بازی مار

فصل سوم: پیاده‌سازی عامل

مدرس:

سید علی کلامی هریس

دانشجوی دکترای حرفه‌ای داروسازی، فعال در حوزه یادگیری عمیق

faradars.org/fvdrlp101

اپیزود و گام

اپیزود (Episode): از شرایط ابتدایی (Initial State) شروع و به شرایط انتهایی (Final State) ختم می‌شود.

گام (Step): هر گام معادل یک بار تصمیم‌گیری و انجام عمل (Action) توسط عامل (Agent) می‌باشد.

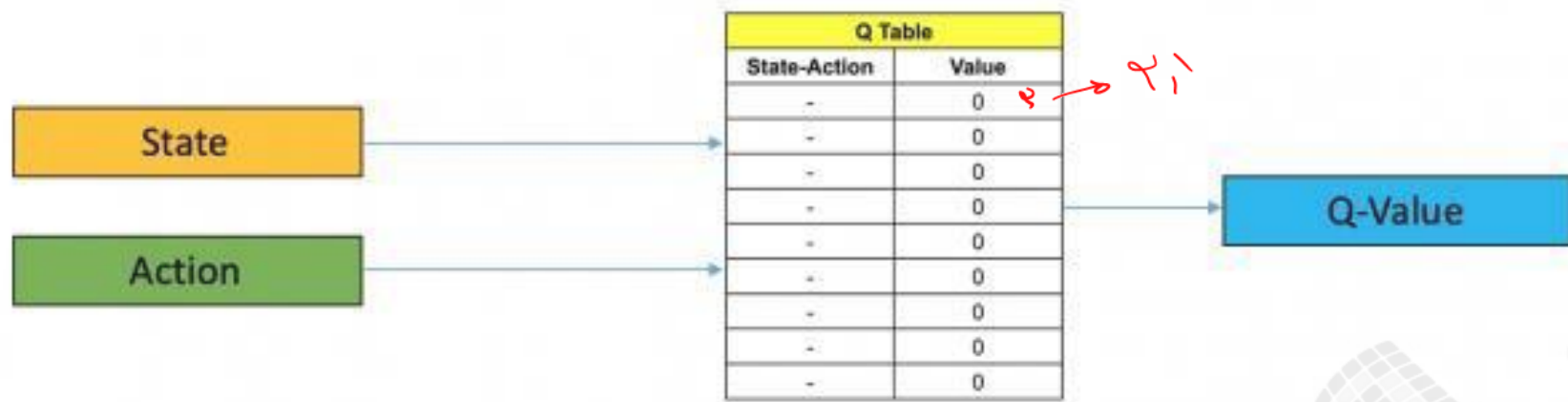
یادگیری Q

در مبحث Q-Learning در انتهای محاسبات به معادله بلمن (Bellman Equation) می‌رسیدیم که اساس آموزش مدل بود:

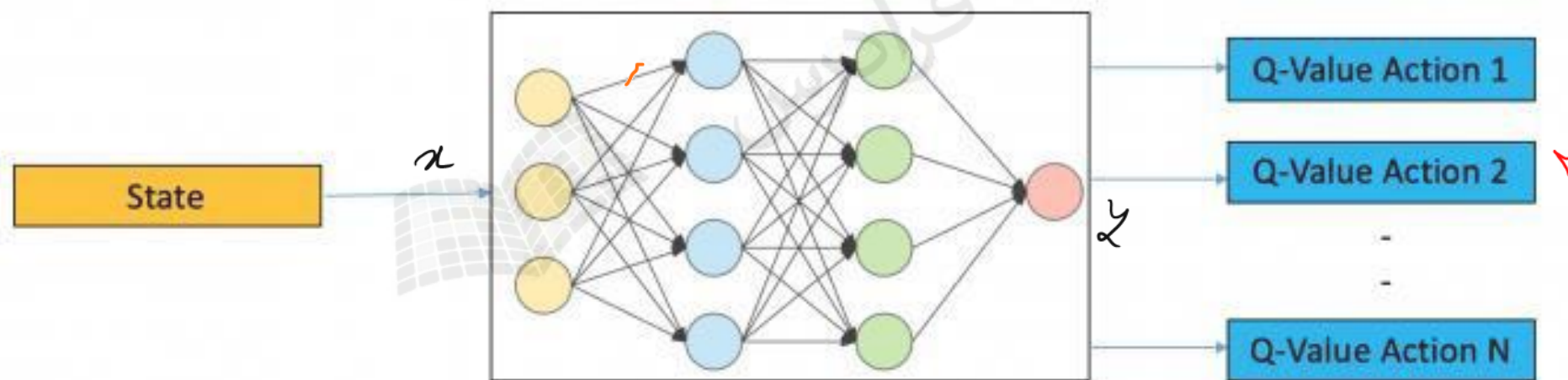
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \cdot \left(r_t + \gamma \cdot \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right)$$

Handwritten annotations in red:

- A bracket above the term $r_t + \gamma \cdot \max_a Q(s_{t+1}, a) - Q(s_t, a_t)$ is labeled "TD".
- An arrow points from the learning rate α to the value "0,1".
- An arrow points from the discount factor γ to the value "0,9".
- An arrow points from the term $Q(s_t, a_t)$ on the right side of the equation to the value "1".



Q Learning

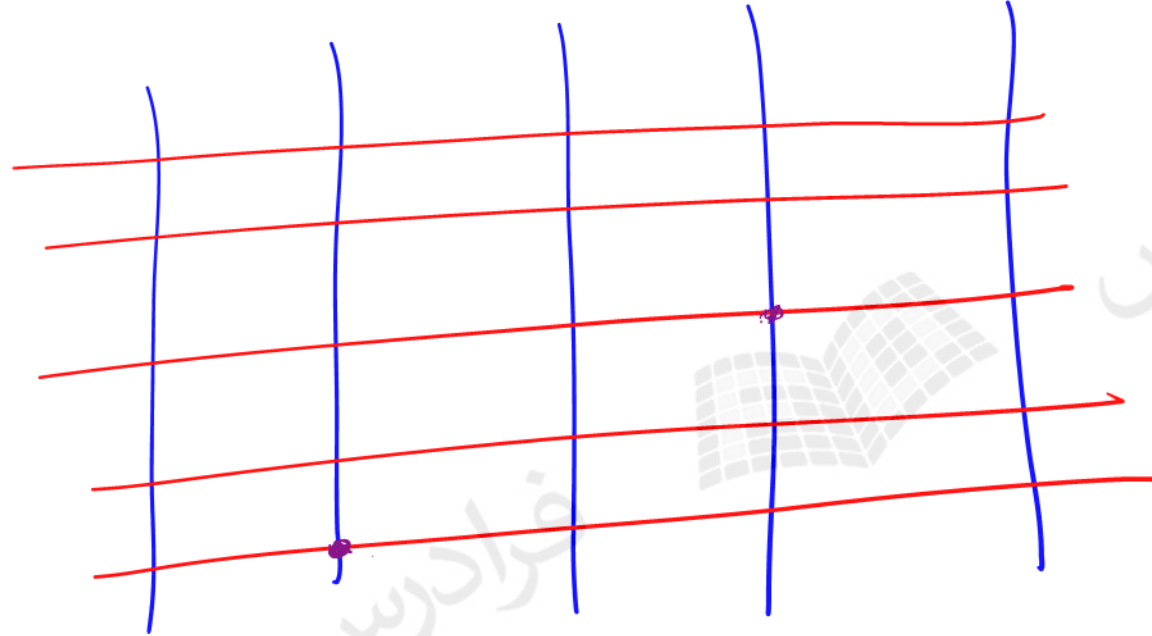


Deep Q Learning

شبکه عصبی

S

A



	+
+	x
x	-
-	.

0

شبکه عمیق Q

برای این مدل اگر یک رشته (Trajectory) به شکل زیر داشته باشیم:

$$s_1, \underline{a_1}, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n$$

پیش‌بینی مدل به شکل زیر خواهد بود:

$$y = r + \gamma \cdot \max_{a'} Q(s', a'; \theta_{t-1})$$

شبکه عمیق Q

تابع زیان (Loss Function) نیز به شکل زیر قابل محاسبه خواهد بود:

$$L(\theta_t) = \mathbb{E}_{s,a,r,s' \sim \rho(\cdot)} \left[(y - Q(s, a; \theta_t))^2 \right]$$

سیاست‌های مورد استفاده

سیاست تصادفی (Random Policy):

$$P(a_t = i) = \frac{1}{n}$$

$$a_t = u\{1, n\}$$

سیاست‌های مورد استفاده

سیاست حریصانه (Greedy Policy):

$$A_m = \underset{a}{\operatorname{argmax}} Q(s_t, a; \theta_t)$$

$$P(a_t = i) = \begin{cases} 1 & i = A_m \\ 0 & \text{otherwise} \end{cases}$$

$$a_t = A_m$$

۱ ۲ -۱ -۲
۰ ۱ ۰ ۰

سیاست‌های مورد استفاده

سیاست حریصانه اپسیلون (Epsilon Greedy Policy):

$$A_m = \operatorname{argmax}_a Q(s_t, a; \theta_t)$$

$$P(a_t = i) = \begin{cases} (1 - \epsilon) + \epsilon/n & i = A_m \\ \epsilon/n & \text{otherwise} \end{cases}$$

$$a_t = \begin{cases} u\{1, n\} \\ A_m \end{cases}$$

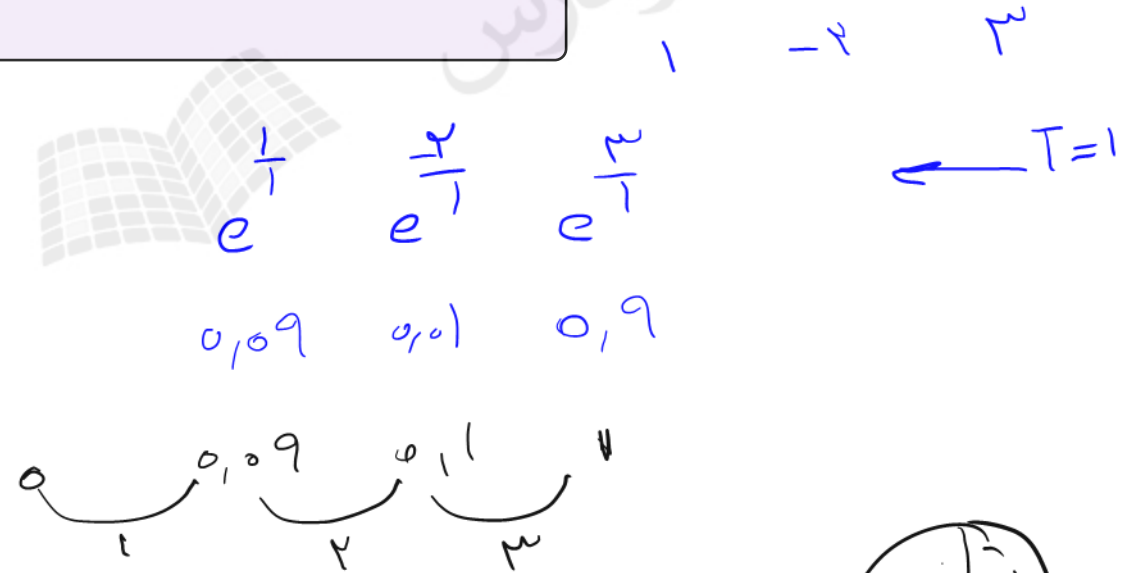
$$U(0,1) < \epsilon \\ \text{otherwise}$$

$\epsilon = 1$	$\epsilon = 0$
$\frac{1}{n}$	1
$\frac{1}{n}$	0

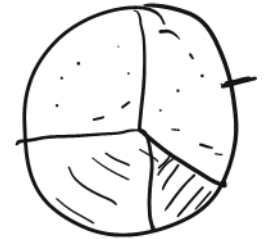
سیاست های مورد استفاده

سیاست بولتزمن (Boltzmann Policy):

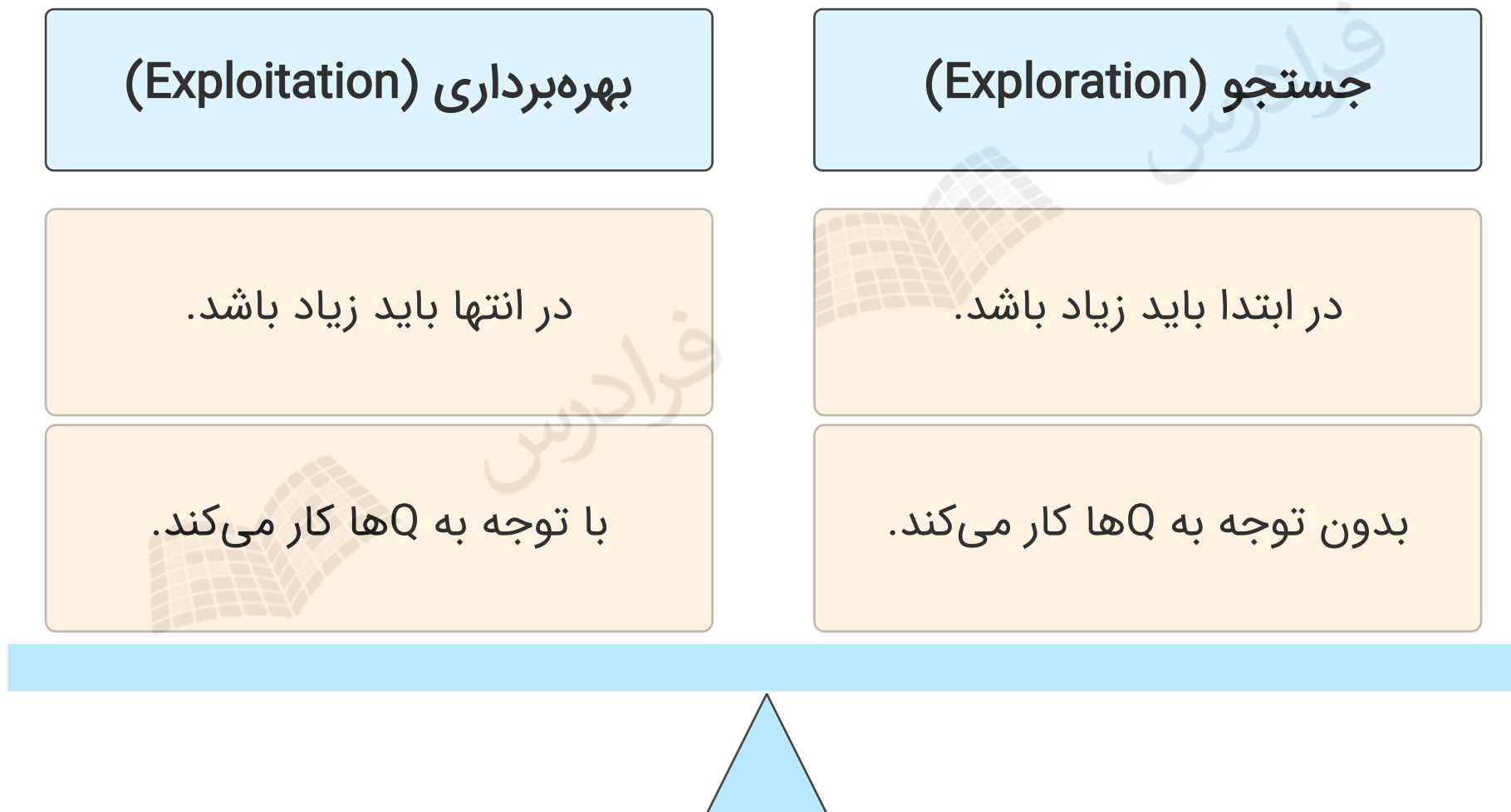
$$P(a_t = i) = \frac{e^{\frac{Q(s_t, i; \theta_t)}{T}}}{\sum_{j=1}^n e^{\frac{Q(s_t, j; \theta_t)}{T}}}$$



$$a_t = \min \left\{ i \mid U(0,1) < \sum_{j=1}^i P(a = j), i \in \{1, 2, \dots, n\} \right\}$$



تعادل بین جستجو و بهره‌برداری

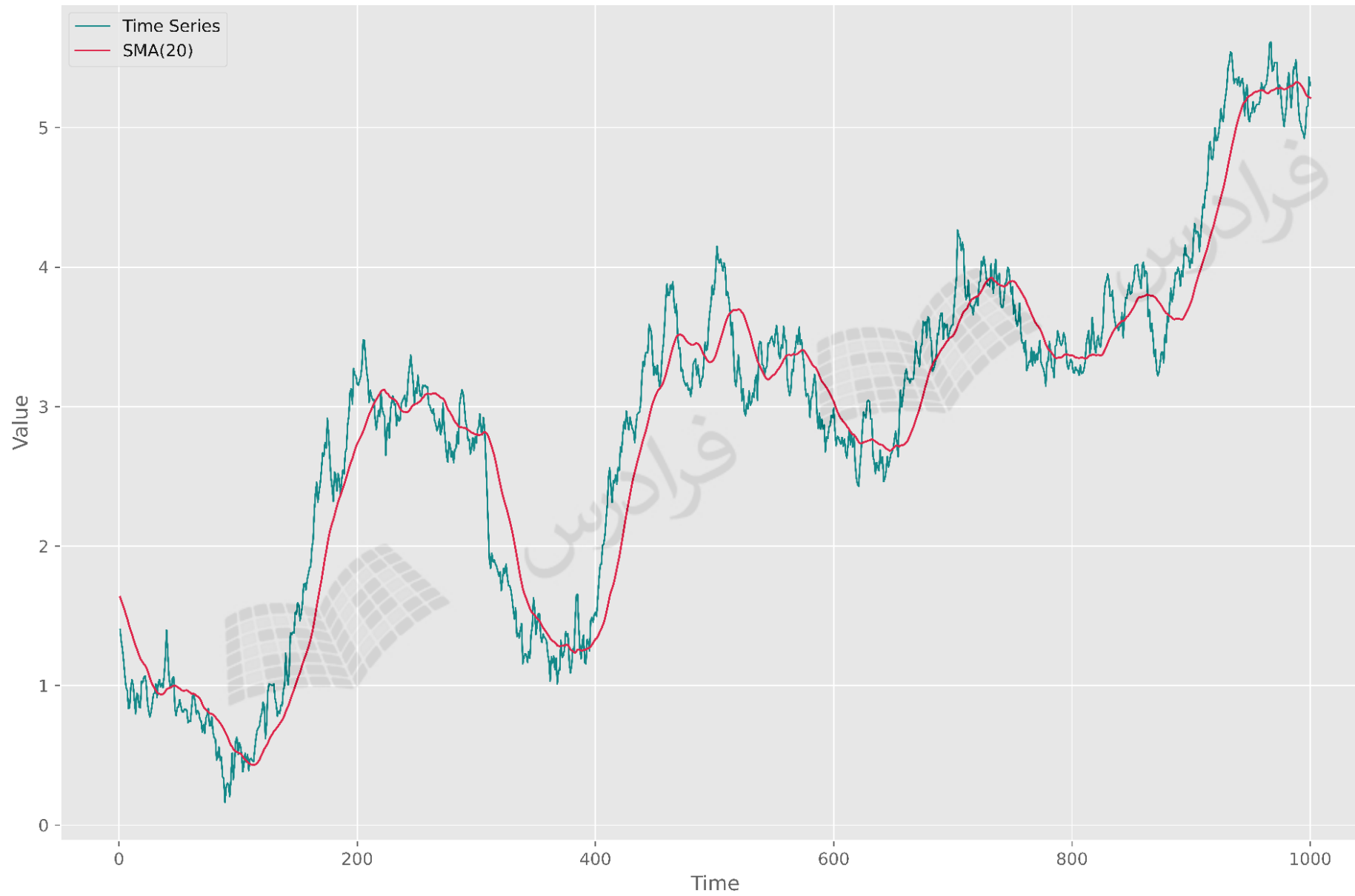


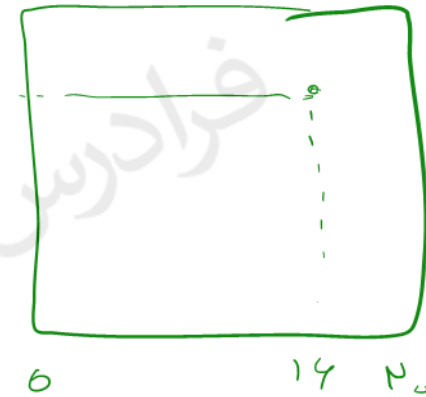
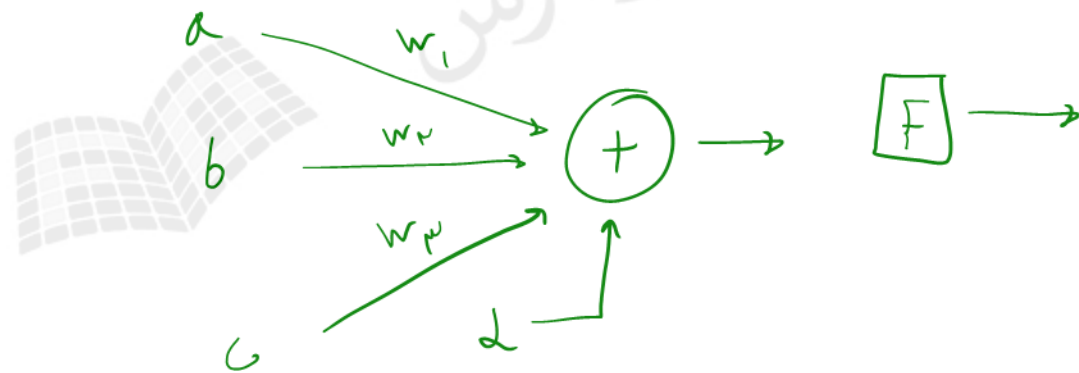
میانگین متحرک چیست؟

میانگین متحرک (Moving Average) یک ابزار بسیار کارآمد می‌باشد که برای نمایش روندها مفید می‌باشد.

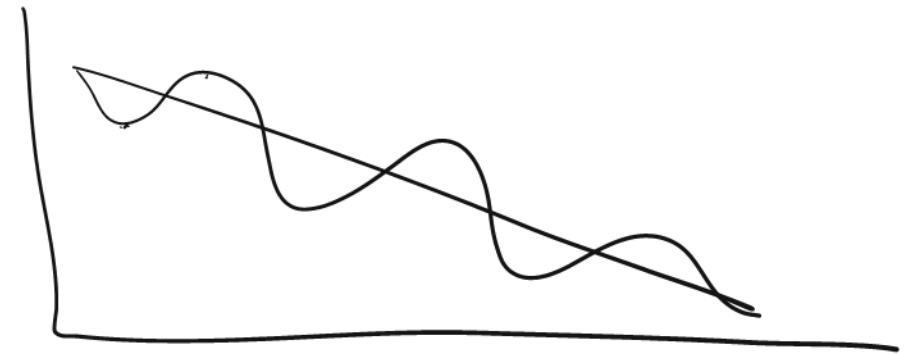
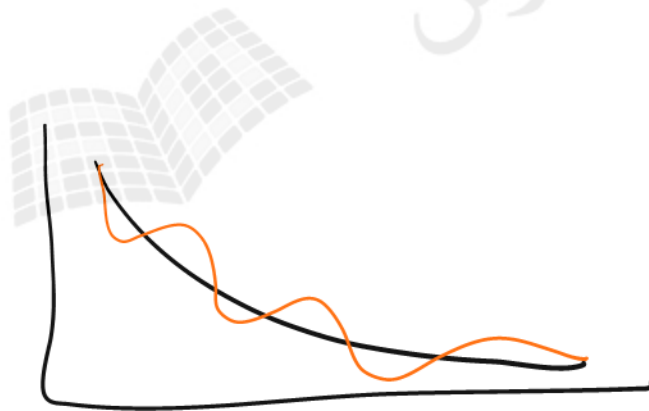
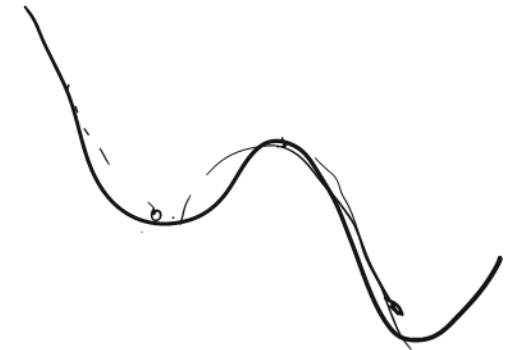
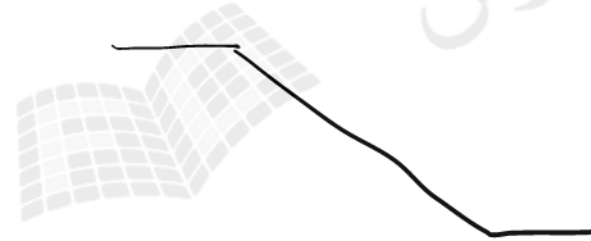
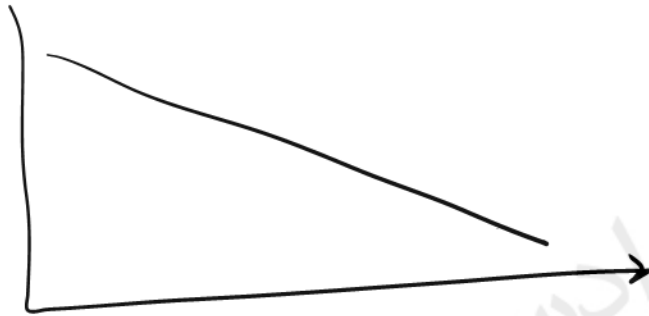
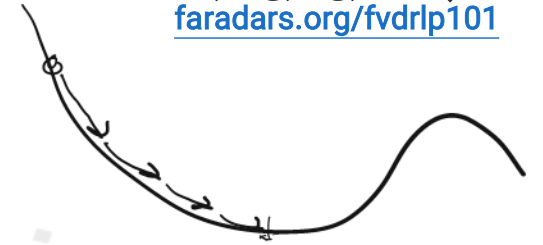
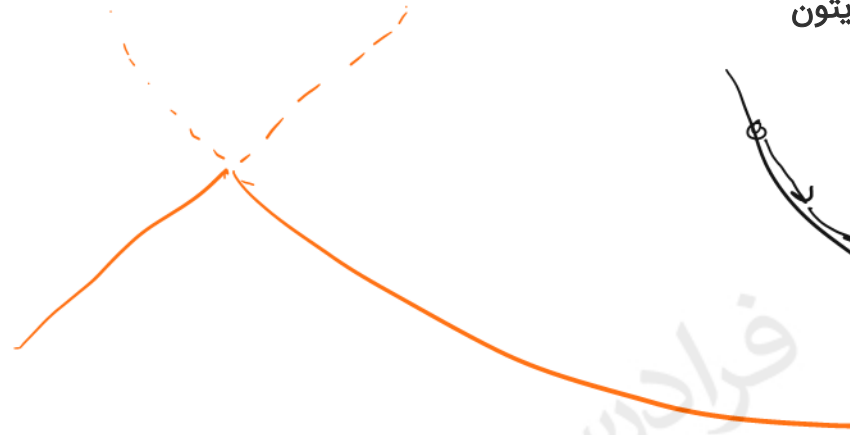
این ابزار یک میانگین با طول پنجره مشخص در طول زمان یا در طول یک آرایه محاسبه می‌کند.

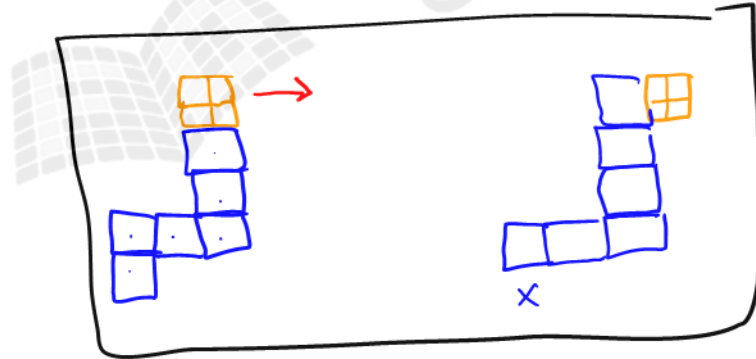
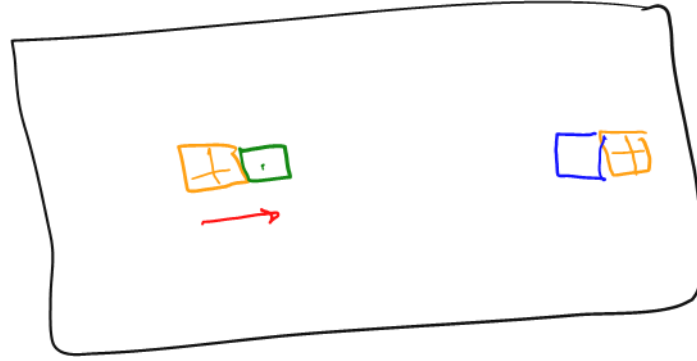
میانگین متحرک، انواع مختلفی دارد که میانگین متحرک ساده (Simple Moving Average یا SMA) ساده‌ترین آن‌ها می‌باشد.



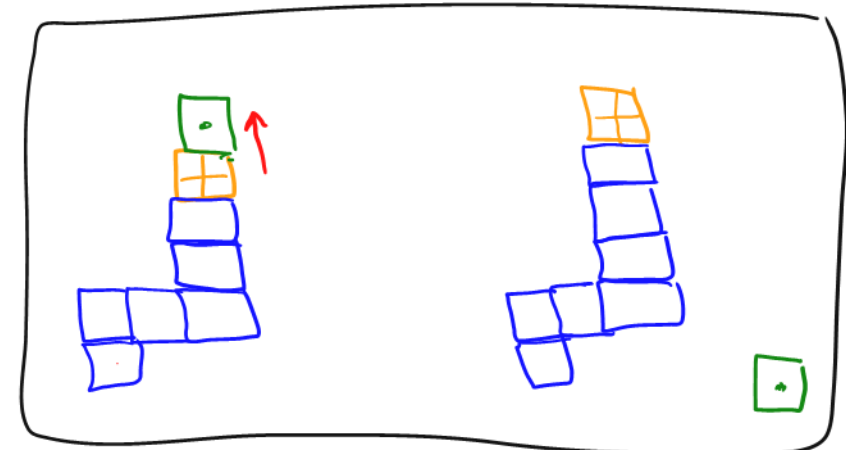


$$\frac{16 - 0}{20 - 0} = \frac{16}{20} = 0.8$$

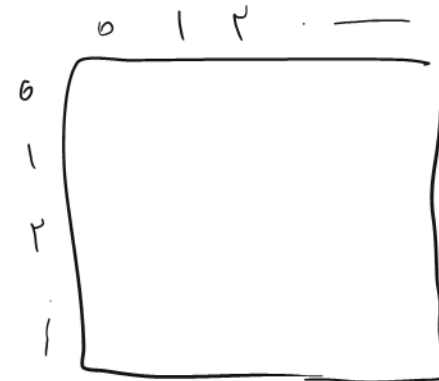
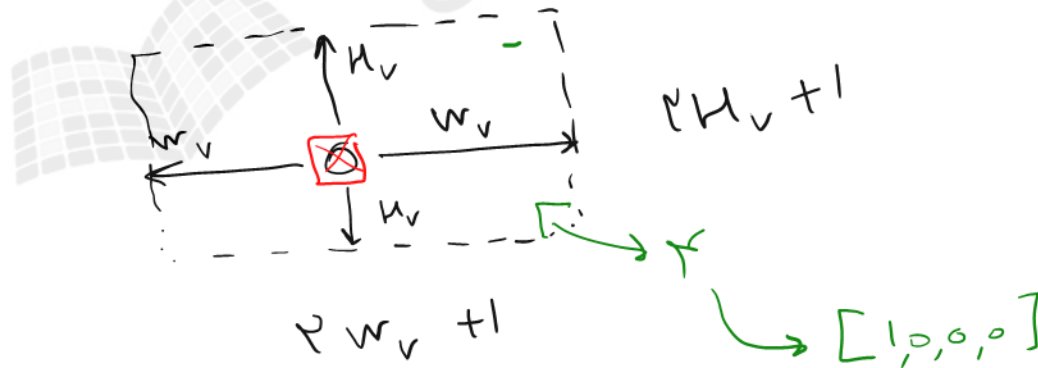
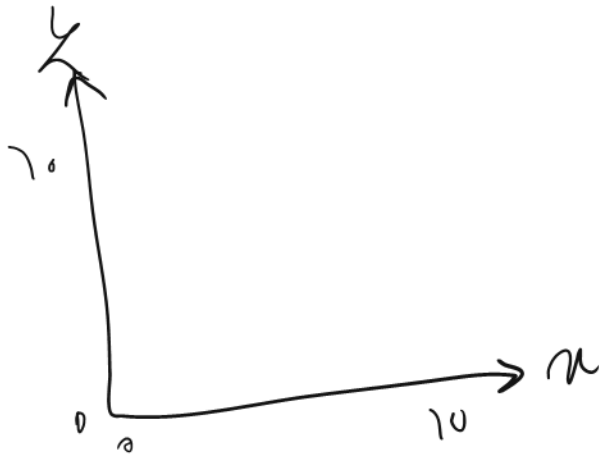




Free



Food



میانگین متحرک چیست؟

اگر یک سری زمانی (Time Series) به شکل زیر داشته باشیم:

$$S = \{S_1, S_2, \dots, S_n\}$$

برای محاسبه میانگین متحرک ساده از رابطه زیر استفاده می‌کنیم:

$$SMA_t = \frac{1}{L} \cdot \sum_{i=1}^L S_{t-i+1}$$

این اسلایدها بر مبنای نکات مطرح شده در فرادرس
آموزش یادگیری تقویتی عمیق در پایتون
پیاده‌سازی بازی مار
تهیه شده است.

برای کسب اطلاعات بیشتر در مورد این آموزش به لینک زیر مراجعه نمایید.

faradars.org/fvdrlp101