

برای تعریف کردن RSTD ابتدا به Reinforcement learning می پردازیم:

وقتی یک لرنینگ برای یک هدف خاصی طراحی نشده باشد اصطلاحاً **unsupervised** است. در یادگیری تقویتی سعی بر این است که از حالت **unsupervised** خارج شویم به طور کامل و بتوانیم **agent** مان را برای هدف خاصی آموزش دهیم.

Classical Conditioning: وقتی کاری که در نظر داریم را انجام داد جایزه‌ای به آن تعلق می‌گیرد و اگر نتوانست انجام دهد تنبیه می‌کنیم.

Operant Conditioning: فعالیت‌های نوروئی یا شبکه‌های نوروئی با دریافت جایزه یا تنبیه اعمال را یاد می‌گیرند. یادگیری تقویتی نه **supervised** است نه **unsupervised**

Unsupervised نیست به این دلیل که به واسطه عملی که **agent** انجام می‌دهد یک فیدبک از محیط می‌گیرد.

Supervised نیست چون به **agent** نمی‌گوییم چی درست یا چی غلط است صرفاً سیگنال‌های جایزه یا تنبیه می‌فرستیم.

یک **agent** یادگیرنده داریم که در یک محیط تعامل انجام می‌دهد. این محیط حالت‌های مختلفی دارد. به عنوان **agent** ابتدا باید محیط را درک کنیم و بدانیم در چه **state** ای وجود داریم. امکان دارد **fully observable** یا **partial observable** باشد. **Agent** تصمیم بگیرد که عملیاتی را در قبال محیط انجام دهد این اکشن محیط را تحت تأثیر قرار میدهد. بر مبنای عملی که انجام داده محیط مجازات می‌کند یا جایزه می‌دهد. تمام حس‌های مثبت و منفی که ادراک می‌کنیم از اعمال و رفتار همان نتیجه فرایند یادگیری تقویتی است. مثلاً سیگنال **reward** در مغز ترشح می‌شود. هر حالت مثبت و منفی که احساس می‌کنیم به خاطر تغییر حالت شیمیایی در مغزمان است. در محیط **Agent** سعی می‌کند جایزه هایش را بیشینه کند.

فیدبک‌ها معمولاً با تاخیر میرسند ← سیستم **credit assignment** برای این که اگر یک پاداش رسید بدانیم برای کدام عمل بوده است.

Distal problem reward: جایزه خیلی فاصله‌ای با عمل ندارد.

- اگر زمان خیلی زیادی بین عمل و جایزه باشد ممکن است **Action** های دیگری هم انجام داده باشیم و ممکن است **credit assignment** سخت باشد که این جایزه یا تنبیه به خاطر کدام عمل بوده که در ازای آن تغییر ایجاد کنیم.
- باید تشخیص دهیم چه بخش‌هایی از شبکه عصبی درگیر این عمل بوده‌اند. که بتوانیم تضعیف یا تقویت انجام دهیم.
- کدام اسپایک‌ها بوده که در ازاش **credit** گرفتیم.

Dopamine: یکی از فرایندهای یادگیری تقویتی شناخته شده در مغز مبتنی بر یک **neuro-modulator** به

اسم **dopamine** است. احساس رضایت در نتیجه دوپامین ایجاد می‌شود.

STDP دوفاز دارد: **LTP**: موقعی که سیناپس‌ها تقویت میشود؛ **LTD**: موقعی که سیناپس‌ها تضعیف شود.

وقتی دوپامین ترشح می‌شود غلظتش می‌تواند نحوه انجام **STDP** را تغییر دهد.

دوپامین می‌تواند یک سری از پترن‌های فایر شدن را reinforce کند و عموماً تأثیر دوپامین روی STDP با تأخیر یک تا دو ثانیه همراه است.

دوپامین همیشه در مغز ترشح می‌شود اگر مقدار آن از یک حدی پایین‌تر بیاید punishment داریم و اگر از یک مقدار بیشتر شود reward.

خاصیتی که دوپامین دارد که بخواهد الگوهای firing که با افزایش دوپامین همراه هستند را تقویت بکند و بقیه تضعیف شوند. بر مبنای یک فرایند سیناپتیکی به اسم synaptic eligibility traces یا synaptic tags است. این فرایند می‌گوید وقتی نورون‌ها فایر می‌شوند تا مدتی بعد از فایر اثر آن در نورون باقی می‌ماند موقعی که دوپامین بعد از تأخیر می‌آید سیناپس‌هایی که Eligibility trace شان بیشتر شده باشد تقویت می‌شوند. اگر دوپامین کم شود، ضعیف می‌شود. به نوعی می‌گوییم STDP را modulate می‌کند.

Reward modulated STDP (RSTDP): دوپامین در فرایند STDP روی دو چیز اثر می‌گذارد: ۱- changing STDP's window ۲- changing STDP's polarity

اگر بخواهیم STDP را به سمتی ببریم که تحت تأثیر دوپامین در مغز چه کاری باید انجام دهیم:

فرض کنیم وضعیت سیناپسها را با دو چیز نشان می‌دهیم:

-وزن سیناپس S

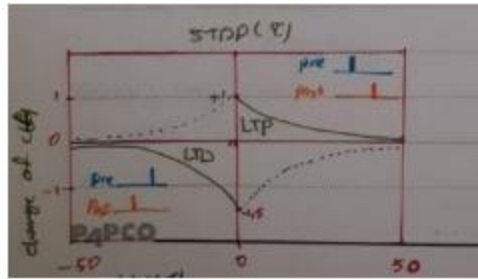
-یک آنزیم ثانویه c که قرار است نقش synaptic tags را اجرا کند.

روابط دینامیک این متغیرها:

$$\frac{dc}{dt} = -\frac{c}{\tau_c} + STDP(\tau)\delta(t - t_{pre/post})$$

$$\frac{ds}{dt} = cd$$

d در واقع میزان دوپامینی است که به صورت خاجی از سلول به نورون رسیده. d اگر منفی باشد دوپامین کم شده و اگر مثبت باشد دوپامین افزایش پیدا کرده. اگر $d > 0$ باشد STDP حالت عادی اجرا می‌شود. اگر $d < 0$ باشد برعکس می‌شود.



اگر نورون اخیرا اسپایک زده باشد $-\frac{c}{\tau_c}$ به صورت exponentially به صفر میل می کند. هر زمان که یکی از نورون های post یا Pre فایر کرده باشد eligibility باید آپدیت شود. $STDP(\tau)\delta(t - t_{pre/post})$ توی لحظاتی که pre یا post فایر کرده باشد آپدیت می کنیم.

هر زمان که سیناپس جایزه یا تنبیه گرفت و c را در d ضرب می کنیم. اگر d مثبت باشد حالت عادی خواهد بود و اگر منفی باشد به anti STDP تبدیل می شود.

پس c میگوید بر مبنای STDP عادی وزن ما باید چقدر دچار LTD یا LTP شود.

یک عامل دوپامین هم هست که میتواند STDP polarity را نگه دارد یا برعکس کند که همان d است.

تفاوت RSTD و STDP در آن است که در RSTD می توانیم عملکرد نورون ها را با سیستم جایزه و تنبیه کنترل کنیم.