

# 相关概念：

blocking:

RPE：主要依赖 error 认知奖励，DA 可算 RPE 生物指标。

{reward-related circuits: DA.

negative reward signals: lateral habenula (LHb): 抑制 DA.

# 强化学习.

① TID learning rule:  $\Rightarrow Q$  似乎与 stratum 活动力相关.

$$Q(S_t, A_t) = E[r(t) + \gamma r(t+1) + \gamma^2 r(t+2) + \dots | S_t, A_t].$$

↓  
期望.  
期望 [ 实际 最终奖励 ]  
+  
未来奖励  
↓ state  
↓ action.

$$\Delta Q(s_t, a_t) = \alpha [r(t) + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)].$$

↓ RPE.      = 学习(实际 + 未来 ↓ ) - 原本预期)

## ②. Exploitation - Exploration.

{ Greedy: best Q.

$\epsilon$ -Greedy: 90% best Q, 10%随机.

Softmax:  $\beta$  { 低  $\beta$ : 随机. (exploration)  
高  $\beta$ : Q 高 (exploitation)

the matching law: 按 reward 比例 做出选择.

③ Future discounting / hyperbolic discounting.

參數：

# Q: 學司率

T：延迟满足能力，越高对未知敏感。

B:  $\beta \uparrow$ , 越偏好高Q.

# 强化学习四大核心模块

模块	关键词	解决的问题	神经实现
RPE (Reward Prediction Error)	$RPE = r + \gamma Q' - Q$	学习应该发生吗? (是否低估/高估了奖励)	多巴胺系统 (VTA/SNc)
Q值与 $\Delta Q$ 更新	$Q \leftarrow Q + \alpha \cdot RPE$	估值如何修正? (调整期望行为价值)	纹状体、皮层突触权重调整
Softmax 动作选择模型	$P(a) = \frac{e^{\beta Q(a)}}{\sum e^{\beta Q}}$	做什么? (概率性行为选择)	vmPFC、前额叶皮层调控 $\beta$
未来奖励贴现模型	$Exp : \gamma^t \quad / \quad Hyp : \frac{1}{1+kt}$	何时行动? (现在好, 还是等待更好)	前额叶-边缘系统、5-HT调控 $\gamma/k$

圖式 schema.

模块	关键词	神经 / 行为机制	易考点
图式 (Schema)	联想网络、结构知识、贝叶斯先验	海马形成、新皮层迁移、情绪联想	图式作用于记忆巩固 / 快速学习
记忆机制	Tse实验、48h迁移、图式一致性	动物paired-associate学习、人类顺序一致任务	图式能加速皮层整合 / 降低认知负荷 越新越激活海马
决策建模	证据整合模型 (DDM) 、LIP区反应 <i>lateral intraparietal (计算脑区)</i>	累积-跨脑决策机制	LIP神经元编码证据方向；行为建模可预测脑区活动 ↑高估小概率, 低估大概率
行为经济学	期望效用、悖论、前景理论	概率权重、损失厌恶、参考点偏好	阿莱悖论解释 / 框架效应实验 / 凹形效用函数 ↓ <i>Petersburg paradox: maximize</i>