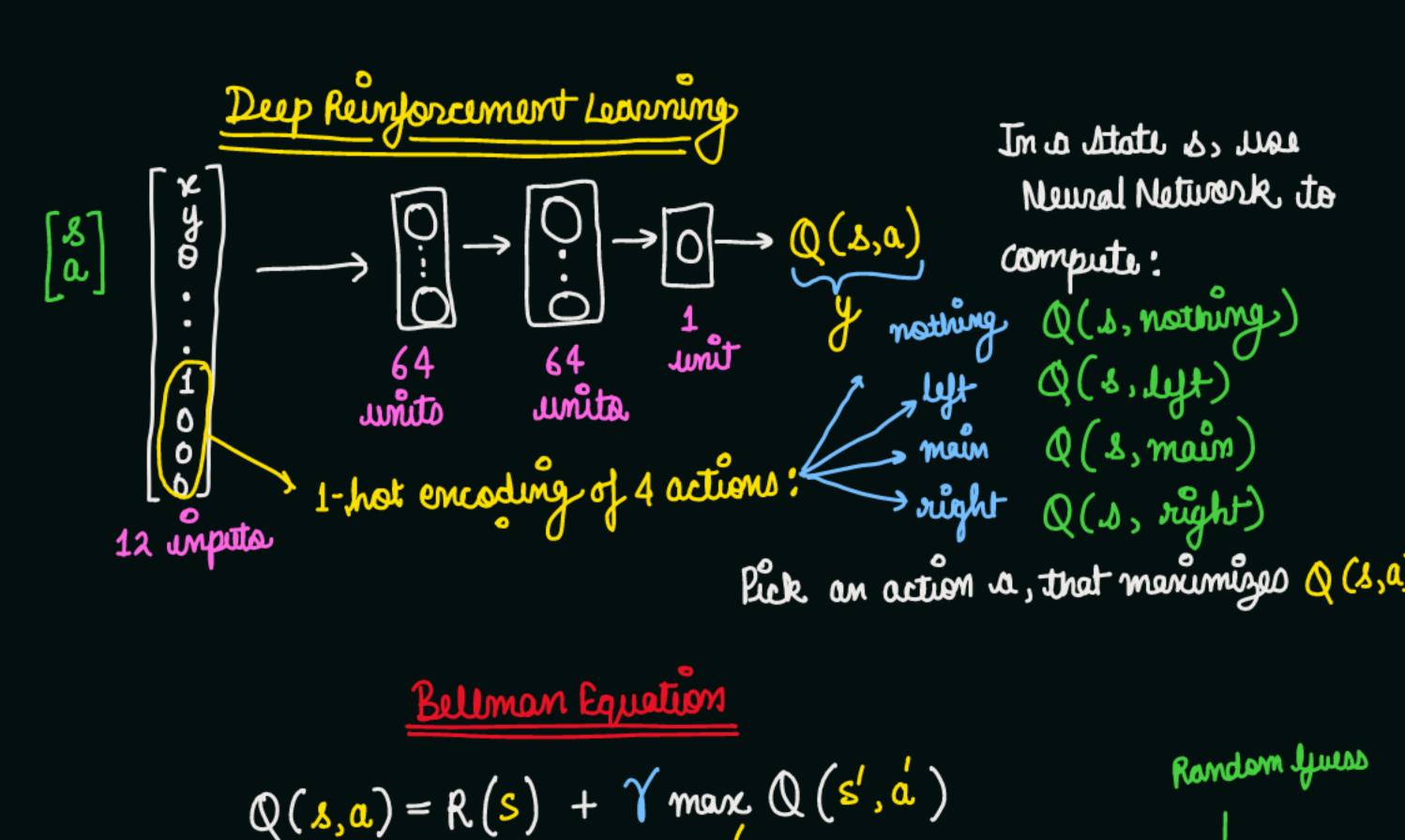
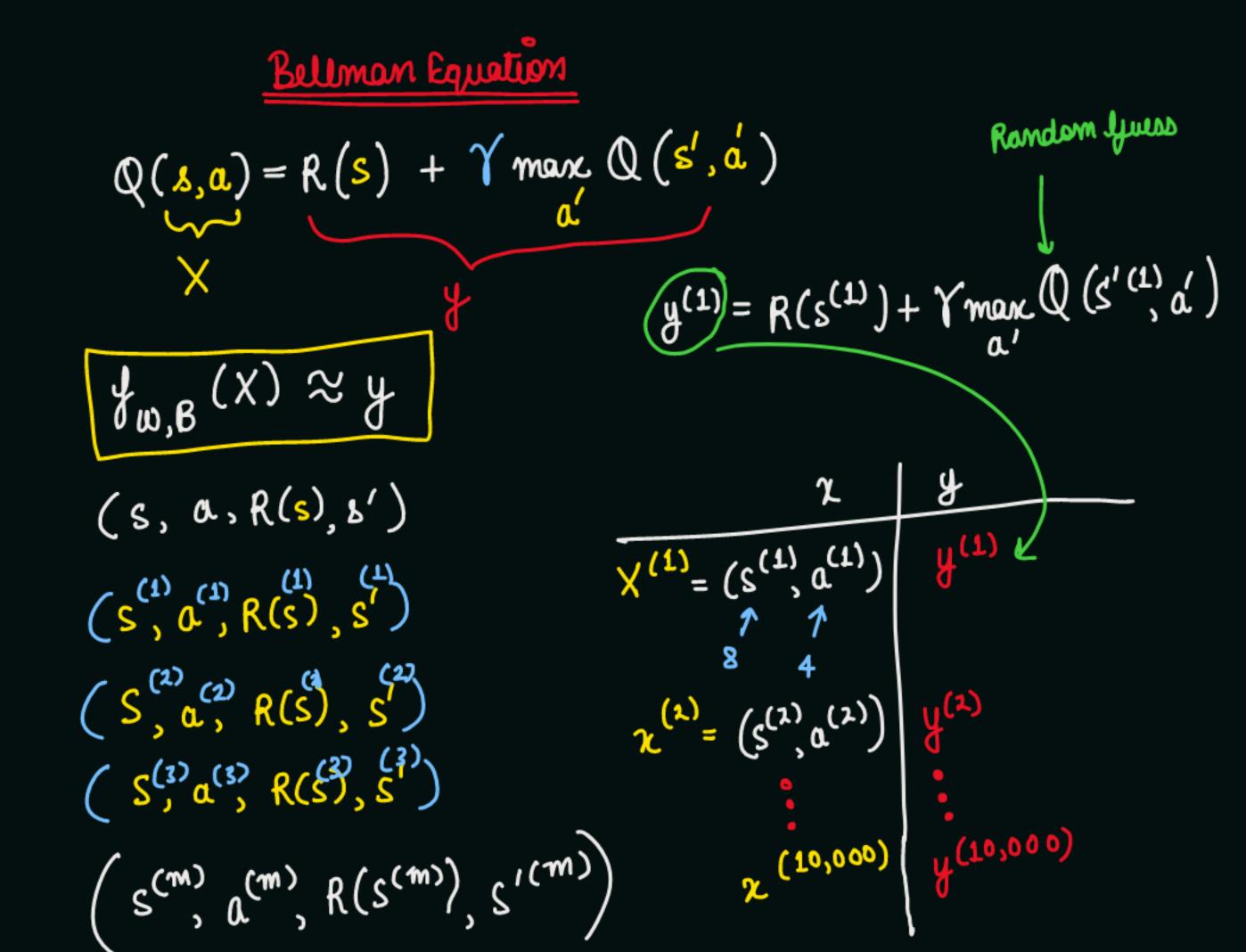
Unsupervised Learning , Recommenders, Reinforcement learning-II Saturday, April 29, 2023 1:06 PM Q(s,a) = Return if you Start in States take action a (once) then behave optimally after that R(s) = reward cof current state &: current state a : current action 8': state you get to after taking action a (2) Factory optimisation (3) Financial (stock) trading (4) Playing games (including video games) $Q(s,a) = R(s) + \gamma \max_{a'} Q(s',a')$ terminal state terminal state $Q(2,\rightarrow) \rightarrow R(2) + 0.5 \text{ mare } (3,a')$ $\rightarrow 0+0.5(25) \Rightarrow 12.5$ (8,a,R(4),4') $Q(4, \leftarrow) = R(4) + 0.5 \text{ maps.}(3, a')$ $(4, \leftarrow, 0, 3)$ b'=3 = 0 + 0.5(25) = 12.5Return = & vell remonde (with discount) $Q(s,a) = R(s) + \gamma \max_{a'} Q(s',a')$ Return = $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots + (until Terminal State)$ > Return from behaving optimally Let Y = 0.9 Return = $0 + 0.9(0) + (0.9^{2}(0) + (0.9)^{3}(100) = 72.9$ (Atate 4 to 1) $\phi(s,a) = R_1 + \Upsilon R_2 + \Upsilon^2 R_3 + \Upsilon^3 R_4 + \cdots$ = $R_1 + \Upsilon R_2 + \Upsilon R_3 + \Upsilon^2 R_4 + ...$ reward of state 8. Return depends on the action you take **太(2)=←** Palicy (T) state policy action ⊼(4) = ← 0.9 0.1 \Rightarrow mapping from state to action, tilling what actions to take in a given state 8. $X(Y) = \sigma$ 0.1 0.9 Expected Return = Average (R1 + YR2 + YZR3 + YR4+...) The goal of RL is to find a policy π , $(a = \pi(A))$ that tells you what action to take in every step to maximize the relation. = $E[R_1 + \Upsilon R_2 + \Upsilon^2 R_3 + \Upsilon^3 R_4 + ...]$ 8 = {1,2,3,4,5,6} Discrete State: 6-stateliner lander 12 3 4 5 6 Continuoua State: x(y)=a Truck/car action a retgailer → do nothing Learn a policy x for given state s, - main thruster pick action, $a = \pi(s)$ - right thruster mouter subsymmen at as as → Crash: -100

mare Q (s,a)

ス(り)=ひ~





m = 100,000,000

(ત્રં,જ)

Batch Learning

new ministration on each iteration

