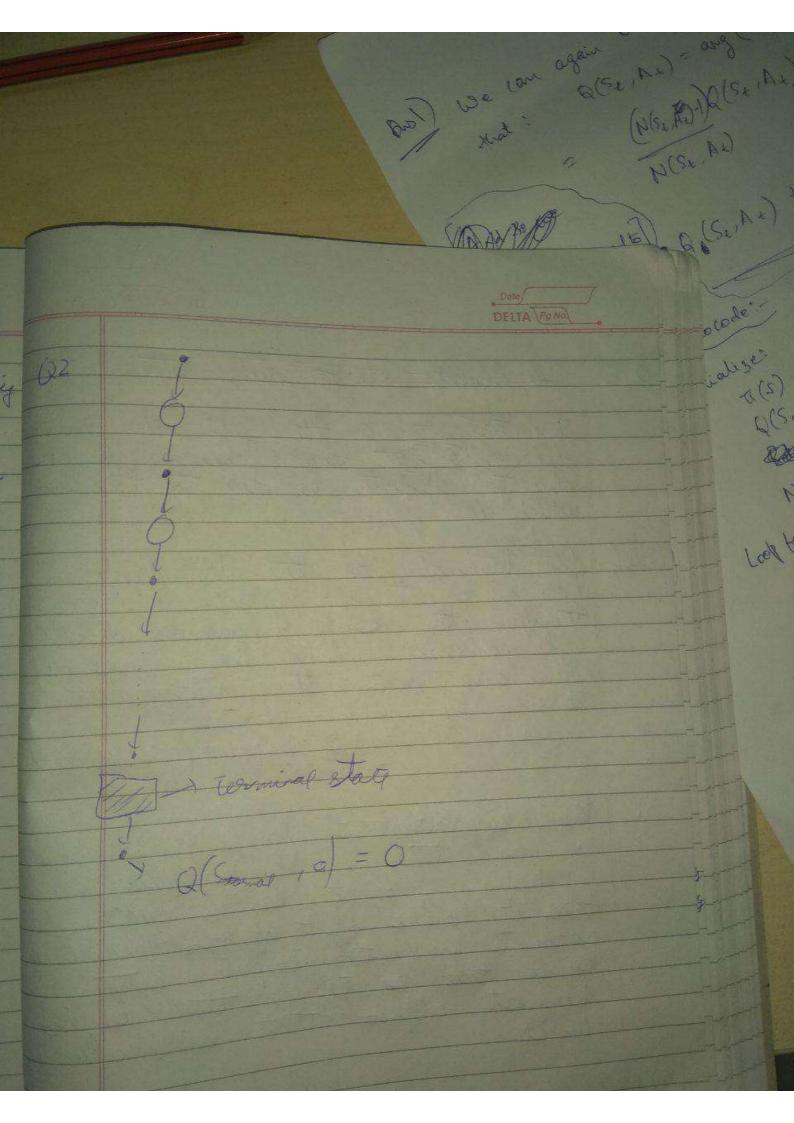
DELTA Pg No Les, as when picking o' to forom policy we will get same texult as when taking argmat. Hence the toes were Q- learning will betome on patity. Q8



Q5 TD defends on bootstrapping, we If we already have statevalues for a model, & then make a minor change to De hours to product of the model, to sing when rober finding a value for this we interest its a value with the a value of initial model, we will see significant gains company to MC since To within the episode to estimate the correct a value to estimate the correct yours of values to estimate the correct value, giving it of the significant increases in according to modified state realized when we are allowed to the significant increases in according to the product of the significant increases in according to the product of the significant increases in according to the significant realized when your good est giving very good est to but & consider the

We went all the way left of stopped all 19: The same here only change only change & We terminated on the left side. For all non-terminal states v.1's are same therefore to one & newends one also can race any effect on these states.

It changed by & (15[x'] + 2+ 10[3]) = \$ (0.150+0-0.5) the olgonoston If longe runber of its petreficial, or it books outsing takes it lowers affect of herent weeks any town the by Thus there is no perfect alpha that will also loore.

Intialize a count (2, a) =0 to £5, @a 61/3) 2 change update of Q(St, to) to Q(St, Az) = Q(St+At) + (G + -Q(S+t)) (A)

Count(St, Az) count(St, At) += 1 Q(8,9) = Ete M(8,0) Pet(t) - Gt Zte Pretty As on In MC, were the episod is already severaled is 6+ is effectively already conditioned on At= a since live know that A = a was the action taken