

PLANNING AND LEARNING

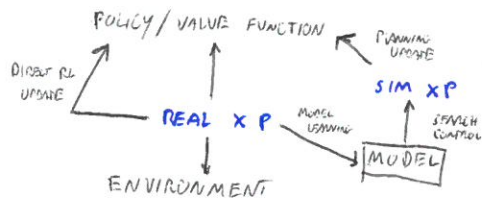
MODEL: SIMULATED EXPERIENCE, SAMPLE MODEL OR DISTRIBUTION MODEL

PLANNING: MODEL $\xrightarrow{\text{PLANNING}}$ NEW/IMPROVED POLICY. STATE-SPACE VS FOM-SPACE PLANNING.

GENERAL STRUCTURE: MODEL $\xrightarrow{\text{SIMULATED EXPERIENCE}}$ BACKUPS $\xrightarrow{\text{VALUES}}$ POLICY

• WHO CARES IF EXPERIENCE IS REAL OR SIMULATED \rightarrow PLANNING \leftrightarrow LEARNING.

DYNA-Q (DYNA AGENTS) • USUALLY SAME RL METHOD • PLANNING | SIM ONLY VISITS PREVIOUSLY OBSERVED STATES AND PREVIOUSLY TAKEN ACTIONS ON S . RANDOMLY.



• DYNA Q+: $Q +$ EXPLOITATION BONUS

• DYNA-AC: ACTION-CRITIC

WHEN ENVIRONMENT IS STOCHASTIC | CHANGES OVER TIME!

• **OPTIMISTIC MODEL:** CHANGES EASILY DISCOVERED BECAUSE EXPLOITABLES RESULT NOW IN WORSE PERFORMANCES.

• **PESSIMISTIC MODEL:** VERY DIFFICULT TO FIND NEW OPPORTUNITIES. • KEEP TRACK OF HOW LONG OF S-A PAIRS HAVE BEEN ATTEMPTED \rightarrow GIVE BONUS POINTS FOR SIMMING THEN

PRIORITIZED SWEEPING:

• UNIFORM SELECTION OF S-A PAIRS TO SIM IS NOT EFFICIENT: WASTEFUL | WASTED BACKUPS.

• **IDEA:** LET'S WORK BACKWARDS FROM STATES WHOSE VALUE HAS CHANGED (IN REAL ENVIRONMENT) AND PROPAGATE TO PREDECESSORS. BACKUP ACTIONS LEADING INTO THAT. KEEP TRACK OF MAGNITUDE OF CHANGES. PRIORITIZE BACKUPS WRT THAT. THRESHOLD.

• VALID IN PRINCIPLE BUT VERY HARD TO IMPLEMENT IN PRACTICE, ALSO RESTRICTED TO FEW CASES

BACKUPS!

• $\{ \text{STATE VALUES, ACTION VALUES} \} \{ \text{OPTIMAL POLICY, ARBITRARY POLICY} \} \{ \text{FULL, SAMPLE} \}$

• 7 OF 8 COMBOS ARE USEFUL ALONGS

• **BRANCHING FACTOR:** NO. OF POSSIBLE NEXT STATES FOR (S, A)

• $b=1$ SAMPLE VS FULL BACKUPS ARE SAME.

• FOR MODERATELY LARGE b , SAMPLE BACKUPS OFFER MORE (S/A) IMPROVEMENT PER UNIT OF COMPUTATION TIME. 'AVALANCH' EFFECT OF EQUAL BACKUPS ESTIMATED TO SUCCESS STATES

TRAJECTORY SAMPLING

HOW TO DISTRIBUTE BACKUPS

EXHAUSTIVE SWEEP (EVERY STATE): VERY EXPENSIVE, UNFOCUSSED, INEFFICIENT

UNIFORM SWEEP: EFFICIENT BUT STILL UNFOCUSSED

ACCORDING TO ON-POLICY DISTRIBUTION: FOLLOWING THE CURRENT POLICY. INTERACT WITH MODEL AND FOLLOW CURRENT π , THEN ON SAMPLING ACTIONS. \rightarrow TRAJECTORY SAMPLING

— EVEN IF WE COULD COMPUTE EXPLICITLY THE DISTRIBUTION, IT WOULD BE INEFFICIENT

+ HELPS SHORT-TERM; ESPECIALLY IF MANY STATES AND LOW b

— MAY HAVE LONG-TERM BECAUSE MORE LIKELY STATES ALREADY HAVE CORRECT VALUES.

HEURISTIC SEARCH

• EXTENDS GREEDY SEARCH BEYOND SINGLE STEP

• IF UNTIL END OF EPISODE \rightarrow WE REEVALUATE OPTIMAL ACTION

• STATE-SPACE SEARCH AS MULTIPLE ONE-STEP BACKUP PIECED TOGETHER.

• DEPTH/SPEED TRADEOFF

• CAN ALSO USE TO DISTRIBUTE BACKUPS \rightarrow USE LIKELY HOOPS OF ACTIONS TO ALLOCATE BACKUPS.

\rightarrow MEMORY AND CPU TIME FOCUS

• GO ON STATES / ACTIONS IMMEDIATELY DEEPER FROM CURRENT SPACE