

RL - APPLICATIONS

TD - GAMMON: $TD(\lambda)$ + NONLINEAR FCN VIA ANN w/ BACKPROP.

SAMUEL'S CHESSERS PLAYER: VERY OLD ('50s), FIRST EFFECTIVE HEURISTIC SEARCH | TD, MINIMAX.

ACROBOT: ROBOT-CONTROL TASK; CONTINUOUS CONTROL VARIABLES, SARSA(λ), LINEAR FCN APPROX, THE GOOD: REDUCING TRACES.

ELEVATOR-DISPATCHING: Q-LEARNING WITH TRUCKS, 1 STEP AHEAD, PPOA-ENCODING CONSTRAINTS.

DYNAMIC CHANNEL ALLOCATION: FOR CELLULAR NETWORKS, SEVERAL ATTEMPTS

JOB-SHOP SCHEDULING: TIME AND RESOURCE CONSTRAINTS \rightarrow OPTIMIZE SCHEDULE \rightarrow **PLAN-SPACE SEARCH!** STATES ARE COMPLETE PLANS

RL PROSPECTS

- I ESTIMATION OF VALUE FUNCTIONS
- II BACKING UP VALUES ALONG ACTUAL OR SIMULATED STATE TRAJECTORIES
- III GENERALIZED POLICY ITERATION. GPI $\pi \rightarrow \hat{V}$

• DIMENSIONS (MAJOR)

- WIDTH OF BACKUPS (SAMPLE VS FULL)
- DEPTH OF BACKUPS (BOOTSTRAPPING)
- FUNCTION APPROXIMATION (TABULAR - LINEAR - NONLINEAR)
- ON-POLICY - OFF-POLICY

• MINOR DIMENSIONS

- EPISODIC | CONTINUING VS DISCOUNT | UNKNOWN REWARDS
- VALUES: ACTION, STATE, AFTERSTATE
- ACTION SELECTION EXPLORATION VS EXPLOITATION
- BACKUPS SYNCHRONOUS | ASYNCHRONOUS
- TRACES NO, ACCUMULATED, REDUCED, DUTCH
- XP REAL, SIMULATED
- BACKUPS: WHAT $S, (S, A)$ SHOULD BE BACKED UP, ACTUALLY ENCOUNTERED OR NOT.
- BACKUPS: TIMING, WHEN ACTION SELECTED OR AFTER
- BACKUPS: MEMORY. PERMANENT RETAINMENT OR ONLY WHEN SELECTING ACTION

• MOVING BEYOND STATE REPRESENTATIONS WITH MARKOV PROPERTY

\rightarrow POMDP, HIDDEN STATES, CAN DO BAYESIAN INFERENCE
- RELIES ON MODEL, COMPUTATIONALLY EXPENSIVE

\rightarrow BETTER REPRESENTATIONS

\rightarrow MODULARITY, HIERARCHY