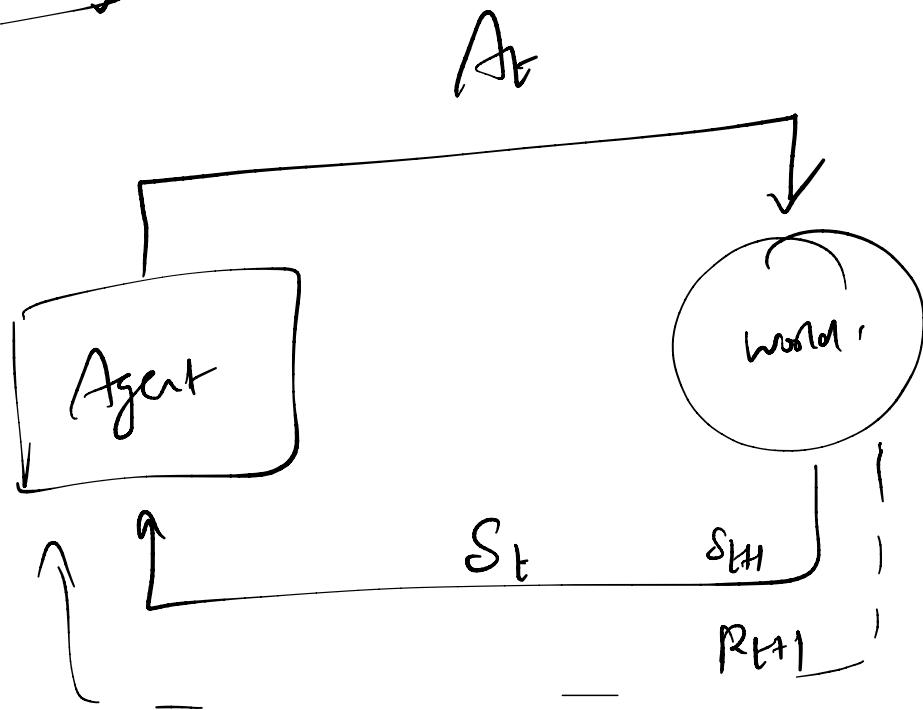


Lecture-13

Reinforcement Learning

→ Learning through interaction.



Maximize the return.

$$R_t + R_{t+1} - \dots$$

→ episodic tasks
 → continuing tasks.

$$R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} - \dots$$

Markov Decision Process: MDP

S, A, T, R, γ

→ State value fn. $V(s)$

→ State-action value fn. $Q(s, a)$

→ Policy $\pi(s)$

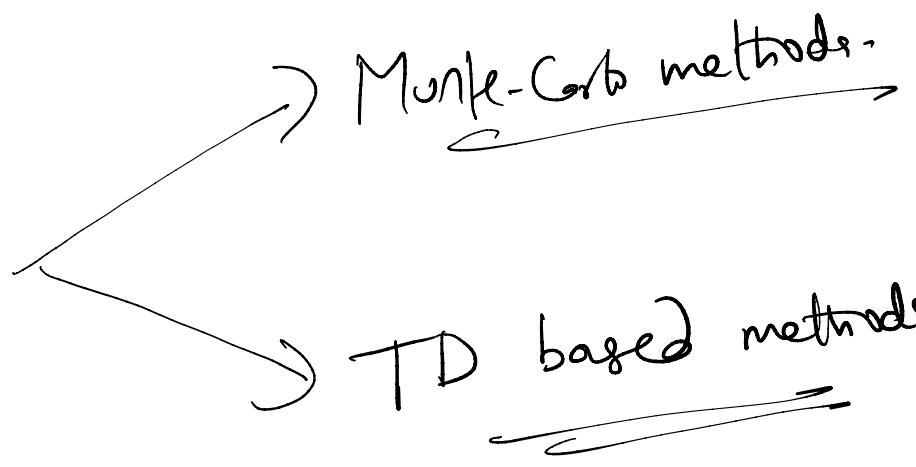
Value-based

Policy gradient
method.

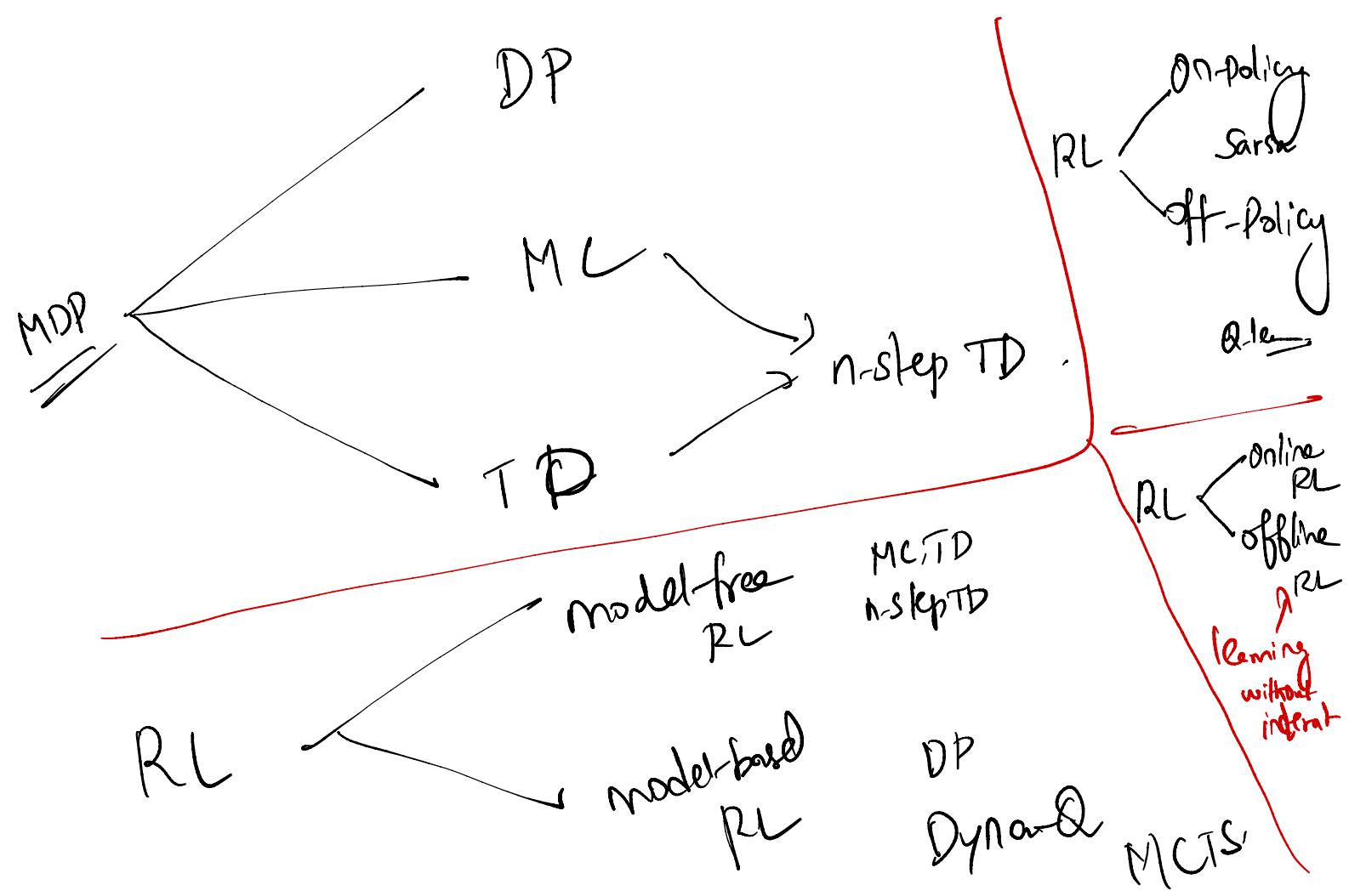
Dynamic Programming

— Bellman equations.

→ Policy iteration.
→ Value iteration.



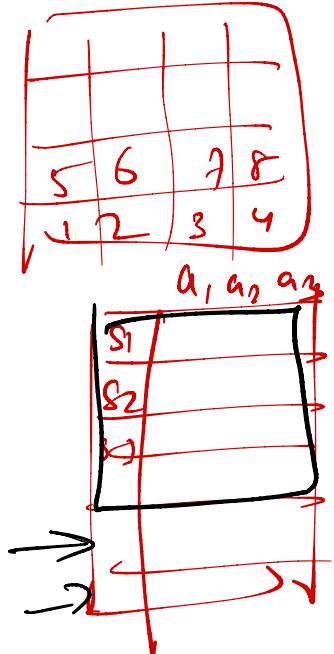
-
- | | |
|---|--|
| MC | TD |
| <ul style="list-style-type: none"> - No bootstrapping, | <ul style="list-style-type: none"> - bootstrapping, |
-



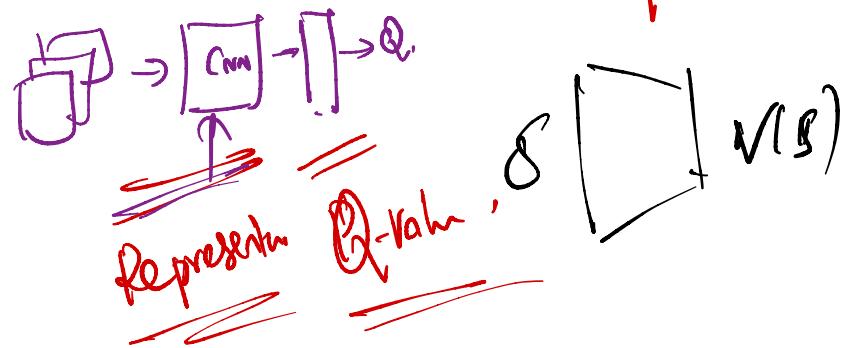
Challenges in Solving RL :-

① Representation learning for RL

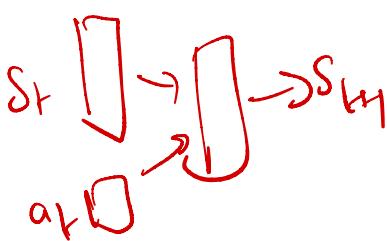
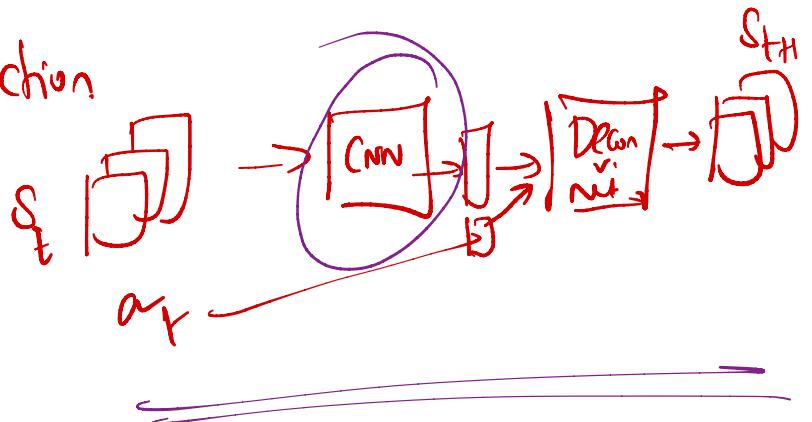
\rightarrow FA \leftarrow Generalizing
Scale



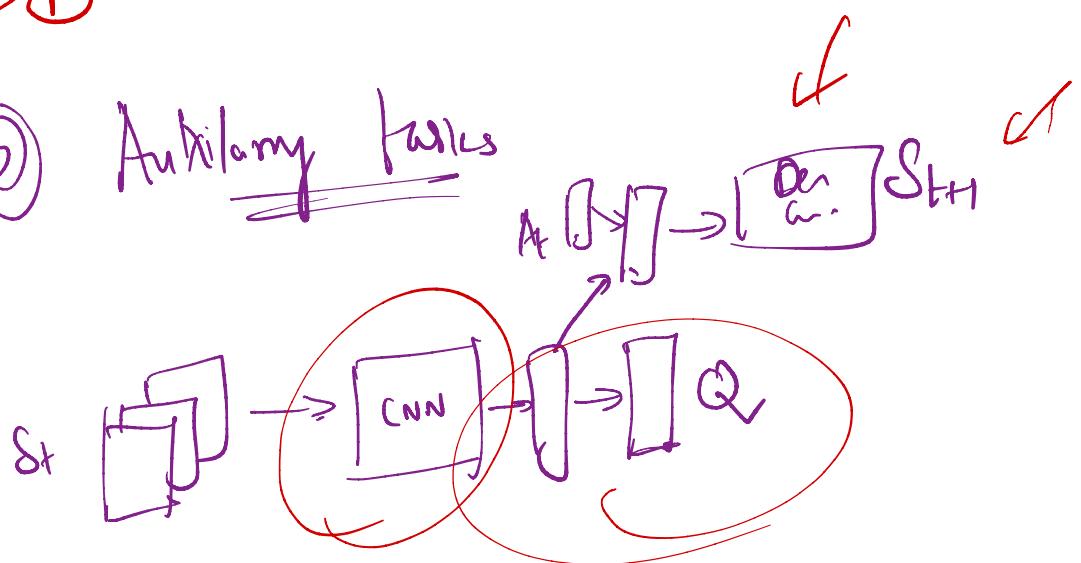
② Pretraining



- forward prediction



③ Antitomy tasks

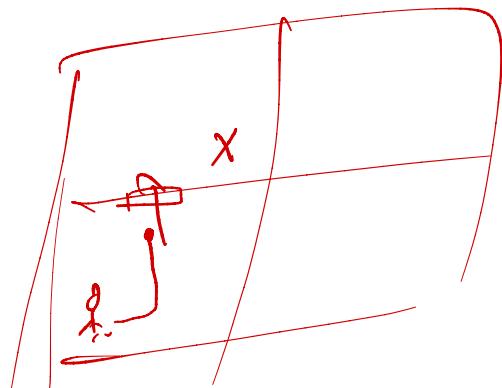


② Partial observability

POMDP
RNNs.

Observations
States

Decision making
under uncertainty



③ State and action abstractions

Hierarchical RL

options
Max Q.

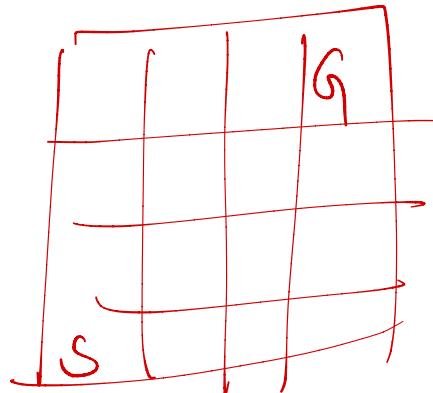
④ Generalization in RL

$$f: x \rightarrow y \quad \boxed{\text{Tr}}$$

Generalize to new states.

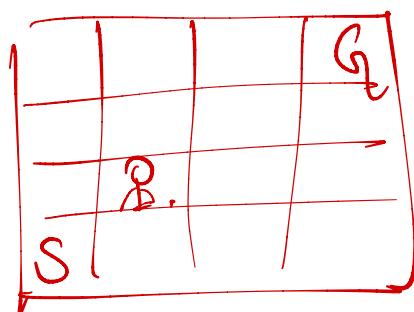
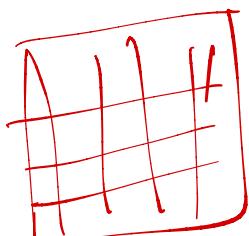
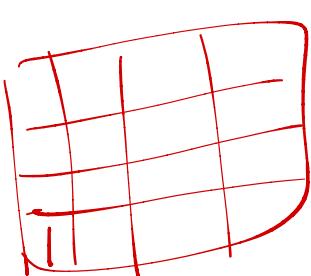
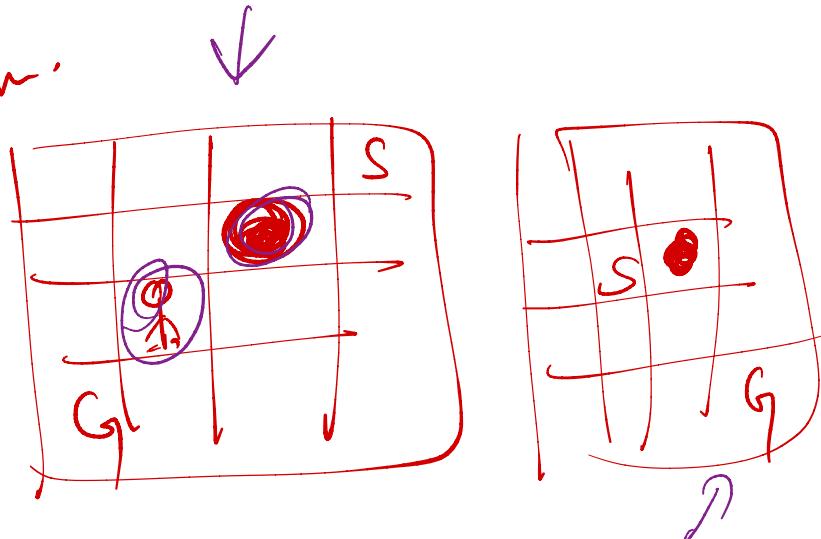
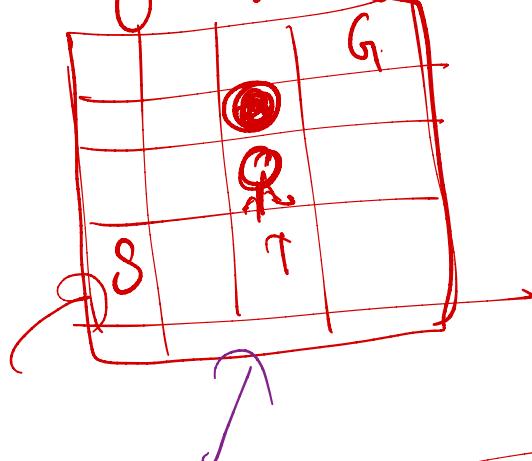
Generalize to new tasks.

Non-stationary env.

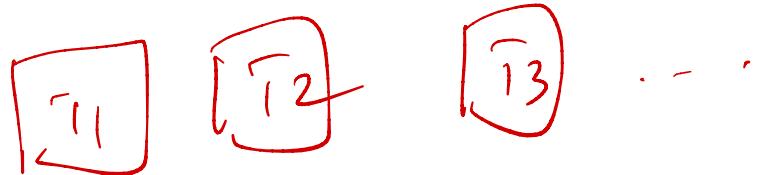


Task: MDP \sim Distribution of MDPs

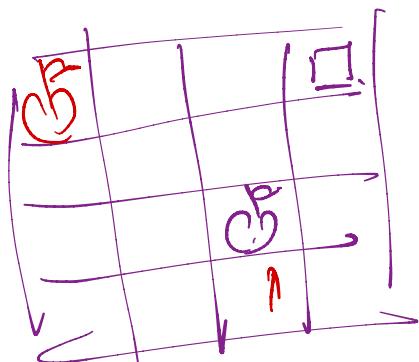
① using right representation.



① → Multi-task RL



→ task interference



so.l.
Apple → +10
box → -10

so.l.
Apple → +10
box → +10

~~so.l.~~ → ~~+10~~

② → lifelong RL / continual RL

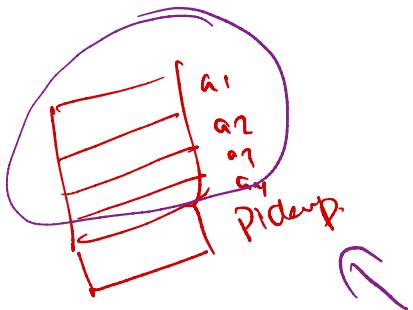
→ Sequence of tasks

① knowledge transfer

Forward transfer

Backward transfer

→ hierarchy / abstraction



→ Self-Supervised learning

Catastrophic forgetting:

→ Continued learning (Udem - next winter)

→ RL Course at Udem (Fall)

→ Deep Learning (Udem → Fall/Winter)

→ RL Course → DeepMind (overlap)

→ Deep RL Course — Sergey Levine
UCB