# Reinforcement learning cont

↳ Actor critic (Both)

**Ex's of RL:** Games (Alpha Go, chess), Robotics,

Finance (stock predictions), Health care, self driving cars (Tesla)
↳ ex of policy based (learn stratigy)
↳ Value based (learn Best moves)

**Algos:** Q-learning, DQN, SARSA, PPO, A3C, RESNFORCE, etc

**challanges:** sample inefficiency → Needs lots of trail and error

Exploration is hard → can get stuck on bad habits

Sparse rewards → Sometimes the agent rarely gets Feedback

Safty → in real world (cars) bad action can be dangeroas

**analogy:** imagine training a dog

• state : Dog sees ball   • action : Dog fetches or ignores

• Reward : Gets reward for Feching   □ over time dog learns Fetch = treat
⟹ reinforced behaviour

**Training** (a training loop is trail and error learning):

1) initialize agent and enviorment
   • agent starts with random policy (no clue what to do)
   • Env is set (eg a game, sim robot)

2) intractive loop
   • agent observes state (s)  • agent chooses action (a) random at first
                                                    (later learned from policy
   • Env responds with: • Reward (r) → (numaric Feedback)  • Next state (s') → new situation

3) learning step (update step)                 Policy parameter (PG) or
   • Agent updates either: Value function (QL/DQL) or   Both (actor critic)
   • goal = maximize expected cumulative reward over time

4) Repeat many episodes: the agent plays episodes (full set of tasks till task ends)
   • with each episode policy improves 102   • Try → Feedback → Adjust → Repeat