# Dynamic Programming (cf Model based)

- ## Optimal Policy

  - policy $\pi$    a distribution over action given states.

    $$\pi(a|s) = P[A_t = a \mid S_t = s]$$

    stationary. time-independent.

    $$P_{ss'}^{\pi} = \sum_{a \in A} \pi(a|s) P_{ss'}^{a}$$
    
    $s \to a$    $s \to a \to s'$

    $$R_s^{\pi} = \sum_{a \in A} \pi(a|s) R_s^a \quad \text{예상 reward}$$

  - ## Optimal Value function

    state-value    $V_*(s) = \max_{\pi} V_{\pi}(s)$     $\longleftarrow \pi \geq \pi'$ if $V_{\pi}(s) \geq V_{\pi'}(s)$, $\forall s$

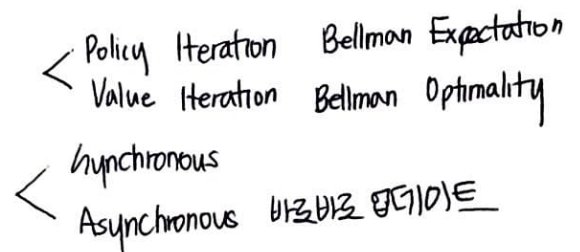    action-value    $q_*(s,a) = \max_{\pi} q_{\pi}(s,a)$.

  - ## Optimal policy can be found by maximising over $q_*(s,a)$

    $$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \arg\max_{a \in A} q_*(s,a) \\ 0 & \end{cases}$$

- ## DP in MDP

  - 대략구조.

    $\left\{\begin{array}{l} \text{Prediction} - \text{Policy Evaluation} \\ \text{Control} \quad - \text{Policy Improve} \quad , \text{Value Iteration} \\ \qquad\qquad\quad \text{-ment} \end{array}\right.$

    $\nearrow$ Policy **Iteration**

    $\left\{\begin{array}{ll} \text{Policy Iteration} & \text{Bellman Expectation} \\ \text{Value Iteration} & \text{Bellman Optimality} \end{array}\right.$

    $\left\{\begin{array}{l} \text{Synchronous} \\ \text{Asynchronous} \quad \text{바로바로 업데이트} \end{array}\right.$

  - ## Policy Evaluation

    : evaluate a given policy $\pi$     $V_1 \to V_2 \to V_3 \to \cdots \longrightarrow V_{\pi}$

    $$\left( \begin{array}{l} V_{k+1}(s) = \sum_{a \in A} \pi(a|s) \left( R_s^a + \gamma \sum_{s \in S} P_{ss'}^a V_k(s') \right) \\ \\ v^{k+1} = R^{\pi} + \gamma P^{\pi} v^k \end{array} \right.$$

    처음에는 random policy 배정, $k$ (횟수) 거듭할수록 optimal policy로.

Control

- Policy Improvement : action 고르기
  ↳ Greedy Policy Improvement

  앞에 $V_{k+1}(s)$ 를 통해 state value 구하기 완료.

  현재 state → 다음 state
                    │
               action 고르기

  $$\pi' = \arg\max_{a \in A} q_\pi(s, a)$$

  $$q_\pi(s, a) = R_s^a + \gamma \sum_{s \in S} \underline{P_{ss'}^a \, V_\pi(s')} \quad \text{뱃룰에 의존}$$

- Value Iteration

  evaluation 에 모든 이동가능한 state 에 대한 state value 구하기 완료

  이중 max 를 취해 greedy하게 value function 을 구하자.

  $$V_{k+1}(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s \in S} \underline{P_{ss'}^a V_k(s')} \right) \quad \text{deterministic}.$$