# CCEN

prof Debbi Prusodi

# INTRODUCTION



PLANT IS UNKNOWN

ACTUATION FUNCTION. ··· PROCESS ··· SENSING FUNCTION.

ACTION $a_t$

MONITORED INFORMATION $H_{t+1}$

MODEL-FREE CONTROLLER (AGENT)

$\overrightarrow{KPI}_{t+1}$ = $\binom{\vec{s}_t}{\vec{a}_t}$

ON-LINE PROCESSING

EHR

R.L. M.D.P.

$f^{NN}$

OFF-LINE PROCESSING

→ SUPERVISED M.L.
→ UNSUPERVISED M.L.

BIG DATA

KNOWLEDGE DATABASE

(DISCRETE TIME)

$s_t$
$r_t$

AGENT

$a_t$

$r_{t+1}$
$s_{t+1}$

ENVIRONMENT

- MACHINE LEARNING TECHNIQUES:
  $\begin{cases} \text{UNSUPERVISED} \longrightarrow \text{PROFILING (OFF-LINE)} \\ \text{SUPERVISED} \longrightarrow \text{NEURAL NETWORKS (ON-LINE)} \\ \text{REINFORCEMENT} \end{cases}$

→ PROFILING: (es/ K-MEANS)

$F_1$ $P_1$ $F_2$ $P_2$ (K=2)
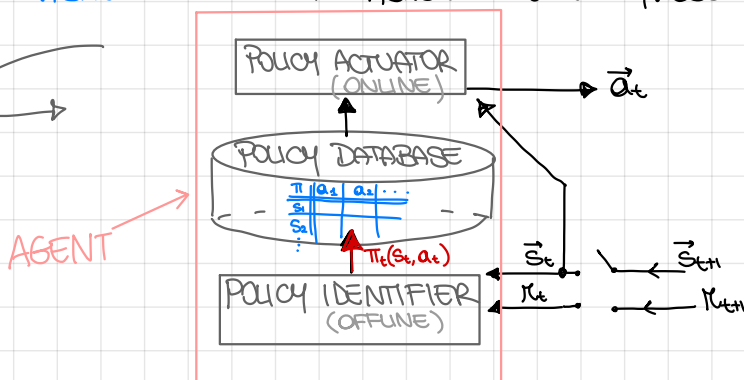
→ NEURAL NETWORKS: GIVEN THE TABLE (INPUT-OUTPUT)

- R.L. = REINFORCEMENT LEARNING → NO KNOWLEDGE OF THE PROCESS IS NEEDED TO CONTROL IT.

- M.D.P. = MARKOV DECISION PROBLEM → KNOWLEDGE NEEDED

STATE → MARKOV PROPERTY: $P\{s_{t+1}|s_t,a_t\} = P\{s_{t+1}|s_t,s_{t-1},\dots;a_t,a_{t-1},\dots\}$

REWARD → MAXIMIZE THE EXPECTED VALUE OF THE LONG TERM RETURN $\max E_{\pi}\left[\sum_{k=0}^{T}\gamma^k \cdot r_{t+k+1}\right]$ $R_t$

↪ YOU CAN ALSO CONSIDER $P\{r_{t+1}|s_{t+1},s_t,a_t\} = P\{r_{t+1}|s_{t+1},s_t,s_{t-1},\dots;a_t,a_{t-1},\dots\}$

ACTION → SELECT THE ACTION TO PERFORM, ACCORDING TO A POLICY $\pi_t(s_t,a_t)$

POLICY ACTUATOR (ONLINE) → $\vec{a}_t$

POLICY DATABASE

$\pi$ | $a_1$ | $a_2$ | ...
$s_1$ |  |  |
$s_2$ |  |  |

AGENT

$\pi_t(s_t,a_t)$

POLICY IDENTIFIER (OFFLINE)

$\vec{s}_t$ $\vec{s}_{t+1}$
$r_t$ $r_{t+1}$

- APPLY THE ACTION AND THEN MONITOR ⇒ INCREASE/DECREASE VALUES ON π-DATABASE

- IF INCREASES IN DIMENSIONALITY ⇒ CAN BE REPLACED BY NEURAL NETWORK.

# MARKOV DECISION PROCESS → YOU KNOW THE PROCESS THROUGH:

- TRANSITION PROBABILITIES: $P_{ss'}^a = P\{S_{t+1}=s' \mid S_t=s, a_t=a\}$

- EXPECTED VALUE OF THE NEXT REWARD: $R_{ss'}^a = E[r_{t+1} \mid S_t=s, a_t=a, S_{t+1}=s']$

  → TARGET: FIND THE OPTIMAL POLICY $\pi^*$ WHICH MAXIMIZE $E_\pi[R_t]$

$$\max \quad E_\pi[R_t] = E_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1}\right] = \sum_{s \in S} \underbrace{E_\pi[R_t \mid S_t=s]}_{= V_\pi^0(s)} P\{S_t=s\}$$

- STATE VALUE FUNCTION FOR POLICY $\pi$:

$$V^\pi(s) = E_\pi[R_t \mid S_t=s] = E_\pi\left[\sum_{k=0}^{\infty} \gamma^k \cdot r_{t+k+1} \mid S_t=s\right]$$

  └→ EXPECTED LONGTERM REWARD WHEN STARTING IN S AND FOLLOWING THE POLICY $\pi$ THEREAFTER.

  └→ **BELLMAN EQUATION:**

$$V^\pi(s) = \sum_{a \in A(s)} \pi(s,a) \sum_{s' \in S} P_{ss'}^a \left[ R_{ss'}^a + \gamma \cdot V^\pi(s') \right] \quad \leadsto \quad \max$$

  $t \to t+1$    $t+1 \to \infty$

  KNOWN     FIXED

  → THE VALUE OF THE STARTING STATE = WEIGHTED SUM OF THE POSSIBLE ARRIVAL STATES
    + THE REWARDS ASSOCIATED TO THE TRANSITION FROM THE START TO THE POSSIBLE ARRIVALS STATES.

  → THE WEIGHT OF EACH POSSIBLE ARRIVAL STATES = PROBABILITY OF ARRIVING AT SUCH STATE.

  → ==HOW TO GET AN OPTIMAL POLICY?== → $\pi^*$ SUCH THAT: $V^{\pi^*}(s) \geq V^\pi(s)$, $\forall s, \forall \pi$.

  (TH. → THERE ALWAYS EXISTS.)

- ACTION-VALUE FUNCTION FOR POLICY $\pi$:

$$Q^\pi(s,a) = E_\pi[R_t \mid S_t=s, a_t=a] \quad \leadsto \quad = \sum_{s' \in S} P_{ss'}^a \left[ R_{ss'}^a + \gamma \cdot \max_{a' \in A(s')} Q^{\pi^*}(s',a') \right]$$

  └→ EXPECTED LONG TERM RETURN WHEN TAKING THE ACTION a IN STATE S AND FOLLOWING THE POLICY $\pi$.

  └→ $V^\pi(s) = \ldots = \sum_{a \in A(s)} \underbrace{E_\pi[R_t \mid S_t=s, a_t=a]}_{= Q^\pi(s,a)} \cdot \underbrace{P\{a_t=a \mid S_t=s\}}_{= \pi} = \sum_{a \in A(s)} Q^\pi(s,a) \cdot \pi(s,a)$

# BACKUP DIAGRAM:

es/ CAN ROBOT ↘ ACTIONS: $A = \{a_1, a_2\}$
STATES: $S = \{s_1, s_2\}$

$t$ · · · · · · · · · · · · · · · ·

$S_1$ · · · · · · · · $S_2$ · · · · ·

$a_1$  $a_2$  $a_1$

$\pi(s_1, a_1) = 0.5$   $\pi(s_1, a_2) = 0.5$   $\pi(s_2, a_1) = 1$

$t+1$ · · ·

$S_1, a_1$  ·  $S_1, a_2$ · · · · · · $S_2, a_1$

$P^{a_1}_{s_1 s_1'} = 0.6$   $P^{a_1}_{s_1 s_2'} = 0.4$   $P^{a_2}_{s_1 s_2'} = 1$   $P^{a_1}_{s_2 s_1'} = 0.5$   $P^{a_1}_{s_2 s_2'} = 0.5$

$S_1'$ · · $S_2'$ · · $S_2'$ · · · · · · $S_1'$ · · · $S_2'$ · · · · ·

$R^{a_1}_{s_1 s_1'} = 2$   $R^{a_1}_{s_1 s_2'} = 1$   $R^{a_2}_{s_1 s_2'} = 3$   $R^{a_1}_{s_2 s_1'} = 2$   $R^{a_1}_{s_2 s_2'} = 4$

AGENT    ENVIRONMENT

$t \to t+1$   $t+1 \to \infty$   $\gamma^1$ BECAUSE II STEP

$$\Rightarrow \left\{ V^\pi(s_1) = 0.5 \cdot 0.6 \left[2 + \gamma \cdot V^\pi(s_1')\right] + 0.5 \cdot 0.6 \left[1 + \gamma \cdot V^\pi(s_2')\right] + 0.5 \cdot 1 \left[3 + \gamma \cdot V^\pi(s_2')\right] \right.$$

$$\Rightarrow \left. V^\pi(s_2) = 1 \cdot 0.5 \left[2 + \gamma \, V^\pi(s_1')\right] + 1 \cdot 0.5 \left[4 + \gamma \cdot V^\pi(s_2')\right] \right.$$

YOU CAN SOLVE IT AS $\left\{ \begin{array}{l} V^\pi(s_1) = V^\pi(s_1') = x \\ V^\pi(s_2) = V^\pi(s_2') = y \end{array} \right.$

$\left\{ \begin{array}{l} x = (0.6 + 0.27x) + (0.2 + 0.18y) + (1.5 + 0.45y) \\ y = (1 + 0.45x) + (2 + 0.45y) \end{array} \right.$   $\Rightarrow \left\{ \begin{array}{l} 0.73x - 0.63y = 2.3 \\ -0.45x + 0.55y = 3 \end{array} \right. \Rightarrow \left\{ \begin{array}{l} x = \underline{26.66} \\ y = \underline{27.27} \end{array} \right.$

$$\Rightarrow \quad E_\pi[R_t] = \underline{26.66} \cdot P\{S_t = s_1\} + \underline{27.77} \cdot P\{S_t = s_2\}$$

THE OPTIMAL WILL BE:

- $$V^{\pi^*}(s) = \max_{a \in A(s)} \sum_{s' \in S} P_{ss'}^a \left[ R_{ss'}^a + \gamma \cdot V^{\pi^*}(s') \right]$$ ← OPTIMAL BELLMAN EQUATION
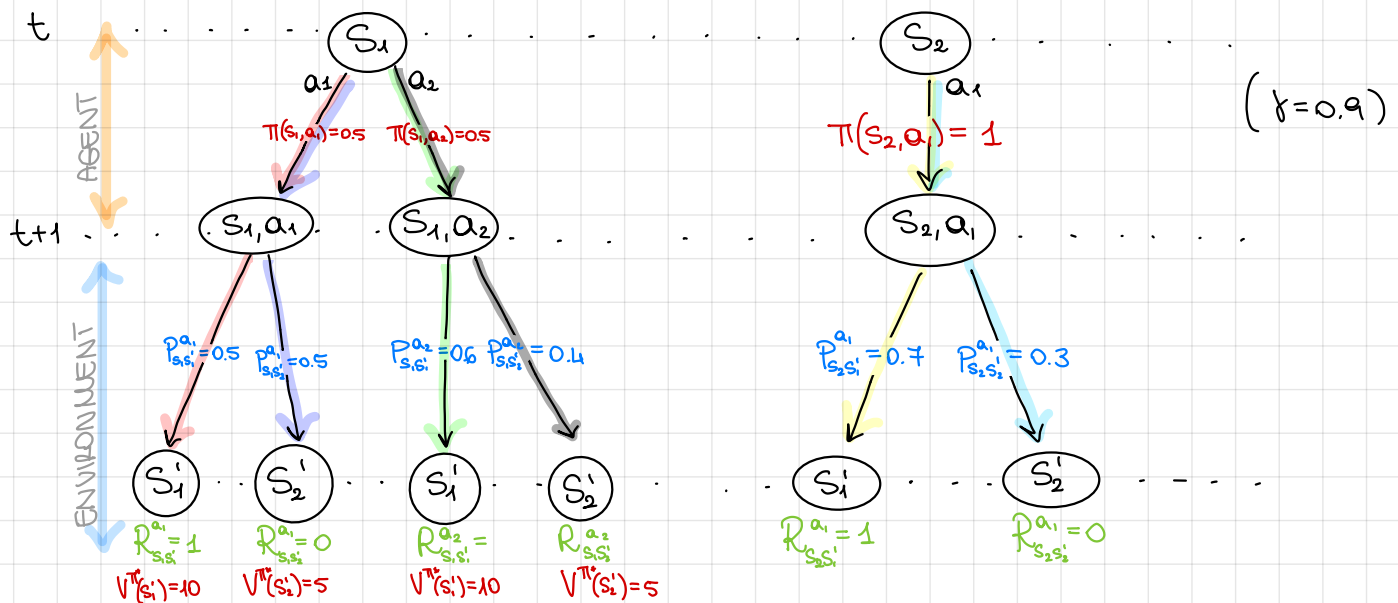
- $$a^*(s) = \arg\max_{a \in A(s)} \sum_{s' \in S} P_{ss'}^a \left[ R_{ss'}^a + \gamma \cdot V^{\pi}(s') \right]$$ ← OPTIMAL ACTION TO PERFORME

- $$Q^{\pi^*}(s,a) = \sum_{s' \in S} P_{ss'}^a \left[ R_{ss'}^a + \gamma \cdot \max_{a' \in A(s')} Q^{\pi^*}(s',a') \right]$$

NB:

When you select the next optimal action it is not optimal just for the next step but it is optimal in a future prospective.

ES/ **BACKUP DIAGRAM:**



$(\gamma = 0.9)$

$(s_1)$
- PATH 🔴 : $1 + 0.9 \cdot 10 = 10$
- PATH 🟣 : $0 + 0.9 \cdot 5 = 4.5$
- PATH 🟢 : $1 + 0.9 \cdot 10 = 10$
- PATH ⚫ : $0 + 0.9 \cdot 5 = 4.5$

- $$V^{\pi^*}(s_1) = \max_{a_1, a_2} \left\{ \underbrace{0.5(1+0.9 \cdot 10) + 0.5(0+0.9 \cdot 5)}_{\text{SELECT } a_1} ; \underbrace{0.6(1+0.9 \cdot 10) + 0.4(0+0.9 \cdot 5)}_{\text{SELECT } a_2} \right\}$$

$$= \max \{ 7.25 ; 7.8 \} \implies 7.8 \implies \underline{a_2} \quad \to \text{SO IF I AM IN } S_1 \text{ I HAVE TO SELECT } a_2.$$

- $$a^*(s_1) = \arg\max_{a_1, a_2} \{ 7.25 ; 7.8 \} = a_2$$

CALL: $\begin{cases} V^{\pi^*}(s_1) = V^{\pi^*}(s_1') = x_1 \\ V^{\pi^*}(s_2) = V^{\pi^*}(s_2') = x_2 \end{cases}$

$$\begin{cases} x_1 = \max \left\{ \overbrace{0.5(1+0.9 \cdot x_1) + 0.5(0+0.9 \cdot x_2)}^{a_1} ; \overbrace{0.6(1+0.9 \cdot x_1) + 0.4(0+0.9 \cdot x_2)}^{a_2} \right\} \quad (s_1) \\ x_2 = 0.7(1+0.9 \cdot x_1) + 0.3(0+0.9 \cdot x_2) \quad (s_2) \end{cases}$$
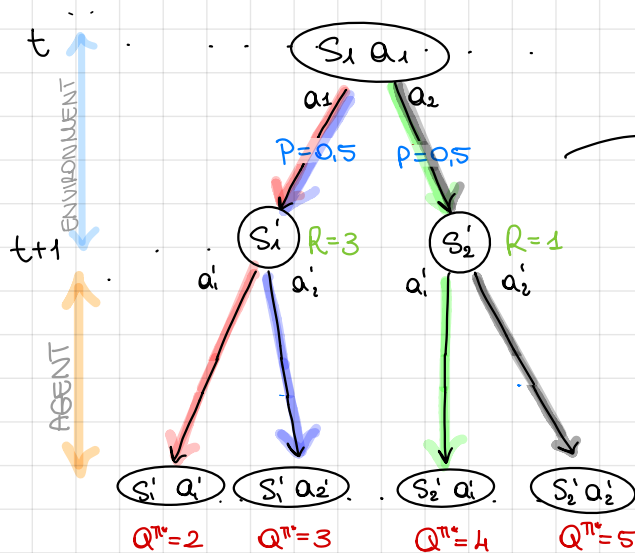
$$\begin{cases} X_1 = \max\{0.45\,X_1 + 0.45\,X_2 + 0.5 \;;\; 0.54\,X_1 + 0.36\,X_2 + 0.6\} \\ X_2 = 0.63\,X_1 + 0.27\,X_2 + 0.7 \end{cases}$$

$\downarrow$ TO SOLVE

- HYP 1 $\rightarrow$ $a_1$ WINS $\Rightarrow \begin{cases} X_1 = 0.45\,X_1 + 0.45\,X_2 + 0.5 \quad = 5.76 \\ X_2 = \qquad '' \qquad\qquad = 5.93 \end{cases}$

- HYP 2 $\rightarrow$ $a_2$ WINS $\Rightarrow \begin{cases} X_1 = 0.54\,X_1 + 0.36\,X_2 + 0.6 \quad = 6.33 \\ X_2 = \qquad '' \qquad\qquad = 6.42 \end{cases}$ HIGHER!

$\Rightarrow$

| $\pi^*$ | $a_1$ | $a_2$ |
|---|---|---|
| $S_1$ | 0 | 1 |
| $S_2$ | 1 | X |

<u>REMARK</u> $\longrightarrow$ I CAN ALSO CONSIDER $Q(s)$ :



$R + \gamma \cdot Q^{\pi}(s,a)$

$$\begin{cases} \text{PATH} \;\bullet : 3 + 0.9 \cdot 10 = 10 \\ \text{PATH} \;\bullet : 3 + 0.9 \cdot 5 = 4.5 \\ \text{PATH} \;\bullet : 1 + 0.9 \cdot 10 = 10 \\ \text{PATH} \;\bullet : 1 + 0.9 \cdot 5 = 4.5 \end{cases}$$

| $Q^{\pi^*}$ | $a_1$ | $a_2$ |
|---|---|---|
| $S_1$ | 24 | 52 |
| $S_2$ | 30 | 18 |

$$\begin{cases} Q^{\pi^*}(S_1, a_1) = Q^{\pi^*}(S_1', a_1') \\ Q^{\pi^*}(S_1, a_2) = Q^{\pi^*}(S_1', a_2') \\ Q^{\pi^*}(S_2, a_1) = Q^{\pi^*}(S_2', a_1') \\ Q^{\pi^*}(S_2, a_2) = Q^{\pi^*}(S_2', a_2') \end{cases}$$

$\rightarrow$ N UNKNOWN
$N = |S| \cdot |A| = 4$

(NB) $\rightarrow$ DIFFERENCE BETWEEN $\longrightarrow$ BELLMAN EQ. :

| $\pi$ | $a_1$ | $a_2$ |
|---|---|---|
| $S_1$ | 0.5 | 0.5 |
| $S_2$ | 1 | / |

OPTIMAL BELLMAN EQ. :

| $\pi^*$ | $a_1$ | $a_2$ |
|---|---|---|
| $S_1$ | 0 | 1 |
| $S_2$ | 1 | / |

# DYNAMIC PROGRAMMING <span>(ch. 4)</span>

↳ HOW TO SOLVE IN A SMARTER WAY → TWO METHODS

## ① POLICY ITERATION → EVALUATION + IMPROVEMENT PROCESS

$$\pi \xrightarrow{\text{EVALUATION}} V^\pi \xrightarrow{\text{IMPROVEMENT}} \pi' \xrightarrow{\text{EVAL}} V^{\pi'} \xrightarrow{\text{IMP}} \pi'' \xrightarrow{\text{EVAL}} \cdots \xrightarrow{\text{IMP}} \pi^*$$

↑ BELLMAN EQ.

$\pi' \geq \pi$
$V^{\pi'}(s) \geq V^\pi(s)$

$\pi'' \geq \pi' \geq \pi$
$\cdots$

WHEN $V^{\bar{\pi}}(s) = V^{\bar{\bar{\pi}}}(s)$

### EVALUATION :

OR

- BY HAND (AS WE ALREADY KNOW)

- COMPUTING $\boxed{V^\pi_{k+1}(s) := \sum_{a \in A(s)} \pi(s,a) \cdot \sum_{s \in S} P^a_{ss'} \left[ R^a_{ss'} + \gamma \cdot V^\pi_k(s) \right]}$, $k = 0, 1, \ldots$

↳ AS $k \to \infty$ → IT CONVERGES TO $V^\pi(s)$

es/ $V^\pi_0(s_1) = V_0(s_1') = 0$, $V^\pi_0(s_2) = V^\pi_0(s_2') = 0$

$\underline{k=0} \rightarrow V^\pi_1(s_1) = 0.55$ ⟿
$V^\pi_1(s_2) = 1.085$

| $V(s_1)$ | $V(s_2)$ |
|---|---|
| $\emptyset$ | $\emptyset$ |
| 0.55 | 1.085 |
| ⋮ | ⋮ |

↓

STOP WHEN $\left| V_{k+1}(s) - V_k(s) \right| < \Delta$  <span>← PREFIXED</span>

### IMPROVEMENT :

$\pi' \geq \pi$ $\begin{cases} \pi'(s, \bar{a}(s)) = 1 \\ \pi'(s, a) = 0 \quad \forall\, a \neq \bar{a}(s), \forall s \end{cases}$

↳ COMPUTE $\boxed{\bar{a}(s) = \arg\max_a Q^\pi(s,a) = \arg\max_{a \in A(s)} \sum P^a_{ss'} \left[ R^a_{ss'} + \gamma \cdot V^\pi(s) \right]}$

es/ $\bar{a}(s_1) = \arg\max_{a_1, a_2} \left\{ Q^\pi(s_1, a_1) ; Q^\pi(s_1, a_2) \right\} = a_2$ ⟹

| $\pi'$ | $a_1$ | $a_2$ |
|---|---|---|
| $s_1$ | 0 | 1 |
| $s_2$ | 1 | / |

→ SINCE $\pi'' = \pi' \Rightarrow \pi' = \pi^*$
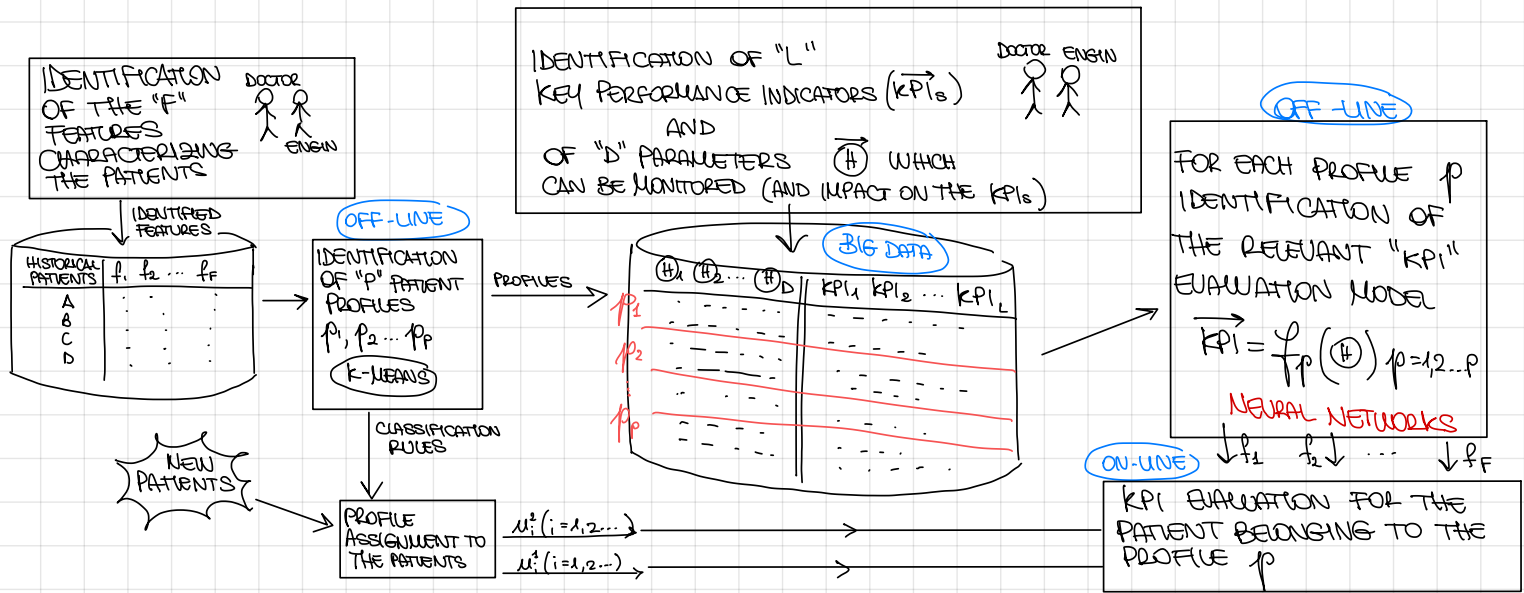
$\bar{a}(s_2) = \cdots = a_1$

## ② VALUE ITERATION

UPDATE THE OPTIMAL BELLMAN EQUATION :

$$\boxed{V^*_{k+1}(s) := \max_{a \in A(s)} \sum_{s \in S} P^a_{ss'} \left[ R^a_{ss'} + \gamma \cdot V^*_k(s') \right]}, \ k = 0, 1, \ldots$$

↳ WHEN $\left| V_{k+1}(s) - V_k(s) \right| < \Delta$ end.

# -CAMS-

IDENTIFICATION OF THE "F" FEATURES CHARACTERIZING THE PATIENTS — DOCTOR ENGIN

↓ IDENTIFIED FEATURES

HISTORICAL PATIENTS | $f_1, f_2 \cdots f_F$
A
B
C
D

IDENTIFICATION OF "P" PATIENT PROFILES $P_1, P_2 \cdots P_P$ (K-MEANS) — OFF-LINE

→ PROFILES → $P_1$ $P_2$ $P_p$

IDENTIFICATION OF "L" KEY PERFORMANCE INDICATORS $(\overrightarrow{KPI}_s)$ AND OF "D" PARAMETERS $\overrightarrow{(H)}$ WHICH CAN BE MONITORED (AND IMPACT ON THE KPIs) — DOCTOR ENGIN

BIG DATA

$H_1$ $H_2$ ... $H_D$ | $KPI_1$ $KPI_2$ ... $KPI_L$

OFF-LINE

FOR EACH PROFILE $P$ IDENTIFICATION OF THE RELEVANT "KPI" EVALUATION MODEL

$$\overrightarrow{KPI} = f_p(\overrightarrow{H}) \quad p=1,2\ldots P$$

NEURAL NETWORKS

↓$f_1$ $f_2$↓ ... ↓$f_F$

NEW PATIENTS

CLASSIFICATION RULES

PROFILE ASSIGNMENT TO THE PATIENTS

$\mu_i^2 (i=1,2\ldots)$
$\mu_i^1 (i=1,2\ldots)$

ON-LINE

KPI EVALUATION FOR THE PATIENT BELONGING TO THE PROFILE $P$

---

## PROFILE 1:

- $f_1$
- $\mu_1^1, \mu_2^1$

$a(\mu_2^1)$ → A — PATIENT $\mu_2^1$ — S → $\overrightarrow{H}(\mu_2^1)$

$a(\mu_1^1)$ → A — PATIENT $\mu_1^1$ — S → $\overrightarrow{H}(\mu_1^1)$

↓ $f_1$

ONLINE
ONLINE $\mu_2^1$

RL CONTROLLERS

$\overrightarrow{KPI}(\mu_2^1)$
$\overrightarrow{KPI}(\mu_1^1)$

KPI EVALUATION FOR THE PATIENT $\mu_1^1$

$$\overrightarrow{KPI}(\mu_1^1) = f_1(\overrightarrow{H}(\mu_1^1))$$

MONITORED INFORMATION

$$\overrightarrow{KPI} = f_i(\overrightarrow{H}) \qquad i = 1,2,\ldots P$$

- TRAINING PHASE → DEDUCE $W_i$ TO REDUCE $f_1$
- OPERATIONAL PHASE → USE $KPI(\mu_i^1) = f_i(\overrightarrow{H})$ [NN]

### INTRODUCE RL:

POLICY ACTUATOR FOR PATIENT $\mu_1^1$ | $\mu_2^1$

$a_t(\mu_2) \to \cdots$
$a_t(\mu_1) \to \cdots$

$\hat{Q}$ FUNCTION EVALUATOR FOR PATIENT $\mu_1^1$ | $\mu_2^1$

$\overrightarrow{KPI}(\mu_2)$
$\overrightarrow{KPI}(\mu_1)$

WITH $\overrightarrow{KPI} = \begin{pmatrix} \vec{s} \\ r \end{pmatrix}$

- $Q_\pi(s_t, a_t) = E_\pi\left[\underbrace{\sum_{k=0}^{\infty} \gamma^k \cdot r_{t+k+1}}_{R_t} \Big| s_t, a_t\right]$ ⤳ $Q^*(s_t, a_t) = \max_\pi Q_\pi(s, a)$

- $r_t = -\|\vec{s_t} - s_{target}\|^2$

RECALL:



$S_t$ $a_t$   $Q(S_t, a_t)$

$S_{t+1}$  $r_{t+1}$

$a_{t+1}$

$S_{t+1}, a_{t+1}$   $Q(S_{t+1}, a_{t+1})$

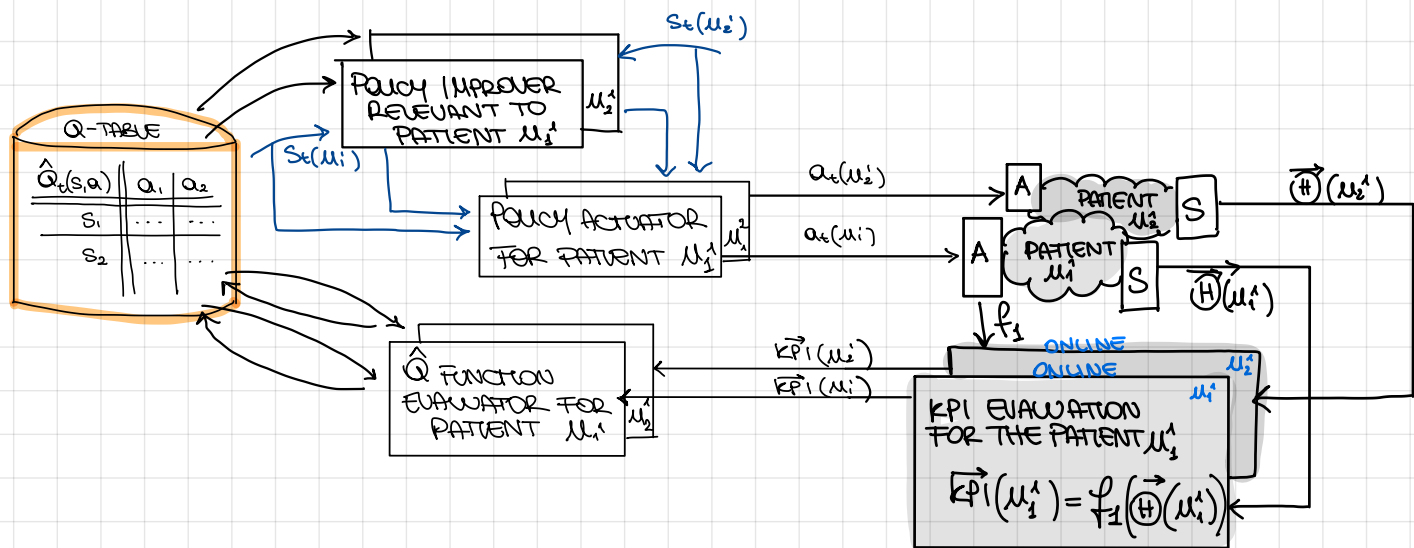$$\alpha_t = \frac{1}{1 + k(s,a)}$$

**Q - LEARNING**

$$\hat{Q}^*(s_t, a_t) = \hat{Q}^*(s_t, a_t) + \alpha_t \left[ r_{t+1} + \gamma \cdot \max_{a' \in A(s_{t+1})} \hat{Q}^*(s_{t+1}, a') - \hat{Q}^*(s_t, a_t) \right]$$

NEW        OLD     STEP SIZE           NEW SAMPLE                OLD

→ TO BUILD THE Q-TABLE



- POLICY IMPROVER COMPUTE THE   $\underset{\text{ACTION}}{\text{PREFERRED}} = \hat{a}_t = \underset{a_t \in A(s_t)}{\arg\max} \; Q^*(s_t, a)$

$$\begin{cases} \pi(s_t, \hat{a}_t) = 1 - \varepsilon + \dfrac{\varepsilon}{|A(s_t)|} \\[2mm] \pi(s_t, a) = \dfrac{\varepsilon}{|A(s_t)|} \end{cases} \quad \text{, WITH} \quad \begin{array}{c} a_t \neq \hat{a}_t \\ a_t \in A(s_t) \end{array} \quad , \quad \varepsilon = \frac{1}{t} \qquad \left( \text{$\varepsilon$ - GREEDY} \right)$$