# Linear MDP

$\phi: S \times A \rightarrow \mathbb{R}^d$

$Q^*$ is linearly realizable    $Q^*(s,a) = \langle \phi(s,a), \theta^* \rangle$ ✓

don't explicitly estimate $P^{t-1}$, dynamic programming for $Q$ directly.

**Recap Tabular   UCB-VI**

$$Q_h^t(s,a) = \min\{H, \ R(s,a) + \sum_{s'} P^{t-1}(s'|s,a) \max_{a'} Q_h^t(s',a') + b^t(s,a)\}$$

bonus
exploration
parameter

For bandits,    $R(s,a) + b^t(s,a)$

For RL (without exp)    $Q = R + \sum_{s'} P(s'|s,a) \max_{a'} Q(s,a')$

Bellman optimality.

UCB-VI update ensures $Q$ is optimistic & nearly Bellman optimal

---

**Linear UCB-VF**    **LSVI-UCB**    $\langle X_i, \beta \rangle$    $Y_i$ ✓

$V_H^t = 0$

$\forall h < H$    $\theta_h^t \leftarrow \underset{\theta}{\arg\min} \sum_{i=1}^{t-1} (\underbrace{\langle \phi(s_h^i, a_h^i), \theta \rangle}_{} - \overbrace{r_h^i - V_{h+1}^t(s_{h+1}^i)}^{})^2 + \lambda \|\theta\|_2^2$

$\Rightarrow Q_h^t(s,a) = \langle \phi(s,a), \theta_h^t \rangle + b_h^{t-}(s,a) \rightarrow \min(H, \dots)$

$V_h^t(s) = \max_{a'} Q_h^t(s,a)$

$\rightarrow$ greedy policy wrt $\theta_h^t$, collect another episode

bonus $b_h^{t-1}(s,a) = \beta \|\phi(s,a)\| \Lambda_{h,t-1}^{-1}$    , $\Lambda_{h,t} = \sum_{i=1}^{t-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^T + \lambda I$

If $H=1$,   LinUCB recovered.

**Complication:**    $E[Y_i] = \langle X_i, \beta^* \rangle$

$E[r_h^i + V_{h+1}^t(s_{h+1}^i)] \rightarrow$ not linear

Least square fit is not guaranteed to be well-specified.

Need more assumptions beyond realizable $Q^*$

**Linear / Low-rank MDP**    if $\exists \phi: S \times A \to \mathbb{R}^d$, $\omega^* \in \mathbb{R}^d$, $\mu: S \to \mathbb{R}^q$
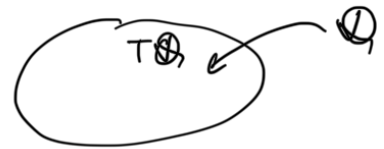
st   $R(s,a) = \langle \phi(s,a), \omega^* \rangle$,   $P(s'|s,a) = \langle \phi(s,a), \mu(s') \rangle$

These assumption imply any $Q$ that satisfies Bellman equation is linear

1) $R(s,a)$ is linear

2) $E_{s' \sim P(\cdot|s,a)}[V(s')]$ is linear

$$= \int_{s'} P(s'|s,a) V(s') = \int_{s'} \langle \phi(s,a), \mu(s') \rangle V(s')$$

$$= \langle \phi(s,a), \underbrace{\int_{s'} \mu(s') V(s')}_{= \bar{\theta}_V} \rangle$$


$T\mathcal{Q}$   $\mathcal{Q}$

**Upper bound:**   $V_p \geq r_s$    LSVI - UCB has $\text{Reg}(T) = \tilde{O}(H^2 \sqrt{d^3 T})$.

   $|S|, |A|$ infinite

   $d$- dim.

**Proof:**   By induction show optimistic regret decomposition for tabular holds.

① Assume $Q^*_{h+1}(s,a) \leq Q^t_{h+1}(s,a)$   $\forall s, a$

$$Q^*_h(s,a) = R(s,a) + \int_{s'} P(s'|s,a) V^*_{h+1}(s')$$

$$= \langle \phi(s,a), \omega^* \rangle + \int_{s'} P(s'|s,a) \max_{a'} Q^*_{h+1}(s',a')$$

$$\leq \langle \phi(s,a), \omega^* \rangle + \int_{s'} P(s'|s,a) \max_{a'} Q^t_{h+1}(s,a) \quad \maltese$$

$$= \langle \phi(s,a), \omega^* \rangle + \langle \phi(s,a), \bar{\theta}^t_h \rangle$$

$$= \langle \phi(s,a), \tilde{\theta}^t_h \rangle \quad \maltese$$

$$= \langle \phi(s,a), \theta^t_h \rangle + \langle \phi(s,a), \theta^t_h - \tilde{\theta}^t_h \rangle$$

$$\underbrace{\leq \|\phi(s,a)\|_{\Lambda^{-1}_{h,t-1}} \|\theta^t_h - \tilde{\theta}^t_h\|_{\Lambda_{h,t-1}}}$$

$$= \langle \phi(s,a), \theta^t_h \rangle + b^t_h(s,a) \quad \checkmark \qquad = \beta_{h,t-1}$$

$$= Q^t_h(s,a) \qquad\qquad \tilde{O}(H^2 d^2) = \beta^2$$

For LS linear fit $\|\hat{\theta}_{LS} - \theta^b\|^2_{\Sigma^{-1}} = \tilde{\sigma}^2 d \rightarrow H^2 d$

$$Y_i = \langle X_i, \beta \rangle + \xi_i$$
$$\downarrow_{\sigma^2}$$

$\sigma \rightarrow H$ + extra $d$ béoz
$\tilde{\theta}^t_h$ is not fixed
as compared to $\theta^{**}$.

②  $Q^t_h(s,a) \le TQ^t_{h-1}(s,a) + \text{conf}^t_h(s,a)$  (nearly Bellman opt)

$$\begin{bmatrix} \text{opt Reg Decomp :} & \text{if } Q \text{ s.t. ① } \theta^* \le \theta^t_h \\[4pt] & \qquad\qquad ② \; Q^t_h \le TQ^t_{h-1} + \text{conf}^t_h \\[6pt] & \text{then Reg} \le \sum_{t,h} \text{conf}^t_h \end{bmatrix}$$

$Q^t_h(s,a) = \langle \phi(s,a), \theta^t_h \rangle + \beta \|\phi(s,a)\|_{\Lambda^{-1}_{h,t-1}}$

$\le \langle \phi(s,a), \tilde{\theta}^t_h \rangle + \langle \phi(s,a), \theta^t_h - \tilde{\theta}^t_h \rangle + \beta \|\phi(s,a)\|_{\Lambda_{h,t-1}}$

$\le \langle \phi(s,a), \tilde{\theta}^t_h \rangle + 2\beta \|\phi(s,a)\|_{\Lambda_{h,t-1}}$

$= TQ^{t-1}_{h-1}(s,a) + \text{conf}^t_h$

Regret $\le \mathbb{E}\left[ \sum_{t,h} \text{conf}^t_h \right] \le 2 \sum_{t,h} \beta \|\phi(s,a)\|_{\Lambda_{h,t-1}} + \tilde{O}(H\sqrt{T})$

↓ using elliptical potential lemma
$\tilde{O}(\beta H \sqrt{dT})$  (from linear bandit)

$= \tilde{O}(H^2 d \sqrt{dT})$

$|S|, |A|$ infinite

**Weaker assumptions.**

$Q^*$ linear realizable

LSVI-UCB
# only requires value
   functions of ↑ form to have
   $\{ \max_a \langle \phi(s,a), \theta \rangle + \beta \|\phi(s,a)\|_{M^{-1}} \}$
   $M \succeq 0, \beta > 0$

TV

**Bellman completeness**    A function class $F: S \times A \to \mathbb{R}$ and MDP $M$ satisfy Bellman completeness if $\forall f \in F$ we have $Tf \in F$.

LSVI-UCB does not work under Bellman completeness.

generalized    $\sigma(\langle \phi(s,a), \theta \rangle)$    included under #
linear                 $\uparrow$ nonlinear
model.

LQR (Linear Quadratic Regulator)    reward-quadratic
                                     transitions-linear

                                    <u>not</u> included under #

Need algorithm that works under weaker Bellman completeness for extension to nonlinear settings in general.