ETH zürich



On the synthesis of Bellman inequalities for data-driven optimal control

Andrea Martinelli • Matilde Gargiani • John Lygeros Automatic Control Laboratory • ETH Zurich

60th IEEE Conference on Decision and Control • Dec. 13-17, 2021



Ingredients:

- A discrete-time system $x^+ = f(x, u, \xi)$ with possibly infinite state & action spaces
- A stage-cost function $\ell : \mathbb{X} \times \mathbb{U} \to \mathbb{R}_+$

Ingredients:

- A discrete-time system $x^+ = f(x, u, \xi)$ with possibly infinite state & action spaces
- A stage-cost function $\ell : \mathbb{X} \times \mathbb{U} \to \mathbb{R}_+$

The discounted ∞ -horizon cost associated to a stationary feedback policy $\pi:\mathbb{X}\to\mathbb{U}$ is

$$m{v}_{\pi}(m{x}) = \mathbb{E}_{\xi}\left[\sum_{k=0}^{\infty} \gamma^k \ell(m{x}_k, \pi(m{x}_k)) \ \Big| \ m{x}_0 = m{x}
ight]$$

Ingredients:

- A discrete-time system $x^+ = f(x, u, \xi)$ with possibly infinite state & action spaces
- A stage-cost function $\ell : \mathbb{X} \times \mathbb{U} \to \mathbb{R}_+$

The discounted ∞ -horizon cost associated to a stationary feedback policy $\pi:\mathbb{X}\to\mathbb{U}$ is

$$m{v}_{\pi}(m{x}) = \mathbb{E}_{\xi}\left[\sum_{k=0}^{\infty} \gamma^k \ell(m{x}_k, \pi(m{x}_k)) \ \Big| \ m{x}_0 = m{x}
ight]$$

Objective: find an optimal policy π^* such that $v_{\pi^*}(x) = \inf_{\pi} v_{\pi}(x) = v^*(x)$



Ingredients:

- A discrete-time system $x^+ = f(x, u, \xi)$ with possibly infinite state & action spaces
- A stage-cost function $\ell : \mathbb{X} \times \mathbb{U} \to \mathbb{R}_+$

The discounted ∞ -horizon cost associated to a stationary feedback policy $\pi:\mathbb{X}\to\mathbb{U}$ is

$$oldsymbol{v}_{\pi}(oldsymbol{x}) = \mathbb{E}_{\xi}\left[\sum_{k=0}^{\infty} \gamma^k \ell(oldsymbol{x}_k, \pi(oldsymbol{x}_k)) \ \Big| \ oldsymbol{x}_0 = oldsymbol{x}
ight]$$

Objective: find an optimal policy π^* such that $v_{\pi^*}(x) = \inf_{\pi} v_{\pi}(x) = v^*(x)$

Classical ADP methods include value iteration, policy iteration and linear programming



Problem setup: the linear programming (LP) approach

The optimal value function admits a recursive definition (the Bellman equation)

$$\mathbf{v}^*(\mathbf{x}) = \underbrace{\inf_{u \in \mathbb{U}} \left\{ \ell(\mathbf{x}, u) + \gamma \mathbb{E}_{\xi} \left[\mathbf{v}^*(f(\mathbf{x}, u, \xi)) \right] \right\}}_{(\mathcal{T}\mathbf{v}^*)(\mathbf{x})}$$

Problem setup: the linear programming (LP) approach

The optimal value function admits a recursive definition (the Bellman equation)

$$\mathbf{v}^*(\mathbf{x}) = \underbrace{\inf_{u \in \mathbb{U}} \left\{ \ell(\mathbf{x}, u) + \gamma \mathbb{E}_{\xi} \left[\mathbf{v}^*(f(\mathbf{x}, u, \xi)) \right] \right\}}_{(\mathcal{T}\mathbf{v}^*)(\mathbf{x})}$$

 $\mathcal T$ is monotone and contractive, hence $\textit{v} \leq \mathcal T\textit{v} \Rightarrow \textit{v} \leq \textit{v}^*$ and therefore

$$\sup_{v\in\mathbb{V}}\int_{\mathbb{X}}v(x)c(dx)$$

s.t. $v(x)\leq (\mathcal{T}v)(x)$ $\forall x,$

Problem setup: the linear programming (LP) approach

The optimal value function admits a recursive definition (the Bellman equation)

$$\mathbf{v}^*(\mathbf{x}) = \underbrace{\inf_{\boldsymbol{u} \in \mathbb{U}} \left\{ \ell(\mathbf{x}, \boldsymbol{u}) + \gamma \mathbb{E}_{\xi} \left[\mathbf{v}^*(f(\mathbf{x}, \boldsymbol{u}, \xi)) \right] \right\}}_{(\mathcal{T}\mathbf{v}^*)(\mathbf{x})}$$

 $\mathcal T$ is monotone and contractive, hence $\textit{v} \leq \mathcal T\textit{v} \Rightarrow \textit{v} \leq \textit{v}^*$ and therefore

$$\sup_{v \in \mathbb{V}} \int_{\mathbb{X}} v(x) c(dx)$$

s.t. $v(x) \leq (\mathcal{T}v)(x) \quad \forall x$

We can relax the constraints by substituting $\ell(x, u) + \gamma \mathbb{E}_{\xi} [v(f(x, u, \xi))] \quad \forall (x, u)$



Efficient methods to solve LPs exist, but several sources of intractability arise

- ▶ v is an optimization variable in the ∞ -dimensional space $\mathbb V$
- \blacktriangleright ∞ number of constraints

$$\sup_{v \in \mathbb{V}} \int_{\mathbb{X}} v(x) c(dx)$$

s.t. $v(x) \le \ell(x, u) + \gamma \mathbb{E}_{\xi} [v(f(x, u, \xi))] \quad \forall (x, u)$

Efficient methods to solve LPs exist, but several sources of intractability arise

- ▶ v is an optimization variable in the ∞ -dimensional space \mathbb{V}
- \blacktriangleright ∞ number of constraints

$$\sup_{v \in \mathbb{V}} \int_{\mathbb{X}} v(x)c(dx)$$

s.t. $v(x) \le \ell(x, u) + \gamma \mathbb{E}_{\xi} [v(f(x, u, \xi))] \quad \forall (x, u)$

We can substitute $\sum_{i} \theta_{i} \phi_{i}(x)$ and sample a finite subset of constraints

$$\mathbf{v}(\mathbf{x}_i) \leq \ell(\mathbf{x}_i, \mathbf{u}_i) + \gamma \mathbb{E}_{\xi} \big[\mathbf{v}(f(\mathbf{x}_i, \mathbf{u}_i, \xi_i)) \big] \quad \forall (\mathbf{x}_i, \mathbf{u}_i) \in \mathcal{D}$$



Efficient methods to solve LPs exist, but several sources of intractability arise

- ▶ v is an optimization variable in the ∞ -dimensional space \mathbb{V}
- \blacktriangleright ∞ number of constraints

$$\sup_{v \in \mathbf{V}} \int_{\mathbb{X}} v(x)c(dx)$$

s.t. $v(x) \le \ell(x, u) + \gamma \mathbb{E}_{\xi} [v(f(x, u, \xi))] \quad \forall (x, u)$
can substitute $\sum_{i} \theta_{i} \phi_{i}(x)$ and sample a finite subset of constraints
 $v(x_{i}) \le \ell(x_{i}, u_{i}) + \gamma \mathbb{E}_{\xi} [v(f(x_{i}, u_{i}, \xi_{i}))] \quad \forall (x_{i}, u_{i}) \in \mathcal{D}$



We

Efficient methods to solve LPs exist, but several sources of intractability arise

- ▶ v is an optimization variable in the ∞ -dimensional space \mathbb{V}
- \blacktriangleright ∞ number of constraints

$$\sup_{v \in \mathbb{V}} \int_{\mathbb{X}} v(x)c(dx)$$

s.t. $v(x) \le \ell(x, u) + \gamma \mathbb{E}_{\xi} [v(f(x, u, \xi))] \quad \forall (x, u)$
We can substitute $\sum_{i} \theta_{i} \phi_{i}(x)$ and sample a finite subset of constraints
 $v(x_{i}) \le \ell(x_{i}, u_{i}) + \gamma \mathbb{E}_{\xi} [v(f(x_{i}, u_{i}, \xi_{i}))] \quad \forall (x_{i}, u_{i}) \in \mathcal{D}$

Model-free/RL framework: construct one constraint for each tuple $\{x_i, u_i, \ell_i, x_i^+\}_{i=1}^T$



Synthesis of Bellman inequalities from data

- Motivation: significant sampling/exploration is required to meet prescribed performance
- How to reconstruct Bellman inequalities from data? When a dataset is sufficiently rich?



Assume for now deterministic linear dynamics and quadratic stage-cost

$$f(x, u, \xi) = Ax + Bu, \quad \ell(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} x \\ u \end{bmatrix}, \quad v \in \mathbb{V}_q = \{v : \mathbb{X} \to \mathbb{R} \mid v(x) = x^{\mathsf{T}} Px\}$$

$$v(x) \leq \ell(x, u) + \gamma \mathbb{E}_{\xi} [v(f(x, u, \xi))]$$

Assume for now deterministic linear dynamics and quadratic stage-cost

$$f(x, u, \xi) = Ax + Bu, \quad \ell(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} x \\ u \end{bmatrix}, \quad v \in \mathbb{V}_q = \{v : \mathbb{X} \to \mathbb{R} \mid v(x) = x^{\mathsf{T}} Px\}$$



Assume for now deterministic linear dynamics and quadratic stage-cost

$$f(x, u, \xi) = Ax + Bu, \quad \ell(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} x \\ u \end{bmatrix}, \quad v \in \mathbb{V}_q = \{v : \mathbb{X} \to \mathbb{R} \mid v(x) = x^{\mathsf{T}} Px\}$$



Assume for now deterministic linear dynamics and quadratic stage-cost

$$f(x, u, \xi) = Ax + Bu, \quad \ell(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} x \\ u \end{bmatrix}, \quad v \in \mathbb{V}_q = \{v : \mathbb{X} \to \mathbb{R} \mid v(x) = x^{\mathsf{T}} Px\}$$

Unknown dynamics: reconstructing H(x, u)

We say that (X, U, X^+) is a dataset of length *T* when

$$X = \begin{bmatrix} x^1 & \cdots & x^T \end{bmatrix}, \quad U = \begin{bmatrix} u^1 & \cdots & u^T \end{bmatrix}$$
 and $X^+ = AX + BU$

Unknown dynamics: reconstructing H(x, u)

We say that (X, U, X^+) is a dataset of length T when

$$X = \begin{bmatrix} x^1 & \cdots & x^T \end{bmatrix}, \quad U = \begin{bmatrix} u^1 & \cdots & u^T \end{bmatrix}$$
 and $X^+ = AX + BU$

Inspired by behavioural theory arguments, Willem's *fundamental lemma* and many related works, we introduce the following result:

Lemma

Consider a dataset (X, U, X^+) of length T with $\begin{bmatrix} X \\ U \end{bmatrix}$ full row-rank. Then, for each $(x, u) \in \mathbb{X} \times \mathbb{U}$ there exists an $\alpha \in \mathbb{R}^T$ such that

 $H(x, u) = (X\alpha)(X\alpha)^{\mathsf{T}} - \gamma(X^+\alpha)(X^+\alpha)^{\mathsf{T}}.$



Unknown stage-cost: reconstructing $\ell(x, u)$

Proposition

Consider a dataset (X, U, X^+) with $\begin{bmatrix} X \\ U \end{bmatrix}$ full row-rank and square. Then, for each $(x, u) \in \mathbb{X} \times \mathbb{U}$ there exists an $\alpha \in \mathbb{R}^{\dim \mathbb{X} \times \dim \mathbb{U}}$ such that

$$\ell(\mathbf{x}, \mathbf{u}) = \alpha^{\mathsf{T}} L_{\mathbf{X}, \mathbf{U}} \alpha, \quad [L_{\mathbf{X}, \mathbf{U}}]_{ij} = \beta \left(\begin{bmatrix} \mathbf{x}_i \\ \mathbf{u}_i \end{bmatrix}, \begin{bmatrix} \mathbf{x}_j \\ \mathbf{u}_i \end{bmatrix} \right),$$

and β is the bilinear form associated to ℓ .

Unknown stage-cost: reconstructing $\ell(x, u)$

Proposition

Consider a dataset (X, U, X^+) with $\begin{bmatrix} X \\ U \end{bmatrix}$ full row-rank and square. Then, for each $(x, u) \in \mathbb{X} \times \mathbb{U}$ there exists an $\alpha \in \mathbb{R}^{\dim \mathbb{X} \times \dim \mathbb{U}}$ such that

$$\ell(\mathbf{x}, \mathbf{u}) = \alpha^{\mathsf{T}} L_{\mathbf{X}, \mathbf{U}} \alpha, \quad [L_{\mathbf{X}, \mathbf{U}}]_{ij} = \beta \left(\begin{bmatrix} \mathbf{x}_i \\ \mathbf{u}_i \end{bmatrix}, \begin{bmatrix} \mathbf{x}_j \\ \mathbf{u}_j \end{bmatrix} \right),$$

and β is the bilinear form associated to ℓ .

• Since
$$\begin{bmatrix} x \\ u \end{bmatrix} = \begin{bmatrix} X \\ U \end{bmatrix} \alpha$$
, we can write $\ell(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} x \\ u \end{bmatrix} = \alpha^{\mathsf{T}} \underbrace{\begin{bmatrix} X \\ U \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} X \\ U \end{bmatrix}}_{L_{x, u}} \alpha$



Unknown stage-cost: reconstructing $\ell(x, u)$

Proposition

Consider a dataset (X, U, X^+) with $\begin{bmatrix} X \\ U \end{bmatrix}$ full row-rank and square. Then, for each $(x, u) \in \mathbb{X} \times \mathbb{U}$ there exists an $\alpha \in \mathbb{R}^{\dim \mathbb{X} \times \dim \mathbb{U}}$ such that

$$\ell(\mathbf{x}, \mathbf{u}) = \alpha^{\mathsf{T}} L_{\mathbf{X}, \mathbf{U}} \alpha, \quad [L_{\mathbf{X}, \mathbf{U}}]_{ij} = \beta \left(\begin{bmatrix} \mathbf{x}_i \\ \mathbf{u}_i \end{bmatrix}, \begin{bmatrix} \mathbf{x}_j \\ \mathbf{u}_j \end{bmatrix} \right),$$

and β is the bilinear form associated to ℓ .

• Since
$$\begin{bmatrix} x \\ u \end{bmatrix} = \begin{bmatrix} X \\ U \end{bmatrix} \alpha$$
, we can write $\ell(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} x \\ u \end{bmatrix} = \alpha^{\mathsf{T}} \underbrace{\begin{bmatrix} X \\ U \end{bmatrix}^{\mathsf{T}} L \begin{bmatrix} X \\ U \end{bmatrix}}_{L_{x, U}} \alpha$
• *L* and *L*_{*X*, *U*} are *congruent* and $\begin{bmatrix} X \\ U \end{bmatrix}$ takes the role of transformation matrix

|U|



2D example



Consider now stochastic linear systems $f(x, u, \xi) = Ax + Bu + \xi$. Estimation with iterative re-initialization can be performed as

$$v(x) \leq \ell(x, u) + \gamma \underbrace{\mathbb{E}_{\xi} \left[v(f(x, u, \xi)) \right]}_{\approx \frac{1}{N} \sum_{k=1}^{N} v(f(x, u, \xi^k))}$$

Consider now stochastic linear systems $f(x, u, \xi) = Ax + Bu + \xi$. Estimation with iterative re-initialization can be performed as

$$\mathbf{v}(\mathbf{x}) \leq \ell(\mathbf{x}, \mathbf{u}) + \gamma \underbrace{\mathbb{E}_{\xi} \left[\mathbf{v}(f(\mathbf{x}, \mathbf{u}, \xi)) \right]}_{\approx \frac{1}{N} \sum_{k=1}^{N} \mathbf{v}(f(\mathbf{x}, \mathbf{u}, \xi^{k}))}$$

Instead, let's write again

$$\mathsf{vec}\big(\underbrace{\mathbb{E}_{\xi}\big[xx^{\mathsf{T}} - \gamma(Ax + Bu + \xi)(Ax + Bu + \xi)^{\mathsf{T}}\big]}_{\mathbb{E}_{\xi}G(x,u,\xi)}\big)^{\mathsf{T}}\mathsf{vec}(P) \leq \ell(x,u)$$



Consider now stochastic linear systems $f(x, u, \xi) = Ax + Bu + \xi$. Estimation with iterative re-initialization can be performed as

$$\mathbf{v}(\mathbf{x}) \leq \ell(\mathbf{x}, \mathbf{u}) + \gamma \underbrace{\mathbb{E}_{\xi} \left[\mathbf{v}(f(\mathbf{x}, \mathbf{u}, \xi)) \right]}_{\approx \frac{1}{N} \sum_{k=1}^{N} \mathbf{v}(f(\mathbf{x}, \mathbf{u}, \xi^{k}))}$$

Instead, let's write again

$$\operatorname{vec}\left(\underbrace{\mathbb{E}_{\xi}\left[xx^{\mathsf{T}} - \gamma(Ax + Bu + \xi)(Ax + Bu + \xi)^{\mathsf{T}}\right]}_{=}\right)^{\mathsf{T}}\operatorname{vec}(P) \leq \ell(x, u)$$



Proposition

Consider a linear stochastic system and $v \in \mathbb{V}_q$. Then, given a dataset (X, U, X^+) ,

$$G\left(ar{x},ar{u},ar{\xi}
ight)=ar{x}ar{x}^{\intercal}-\gammaar{x}^{+}ar{x}^{+\intercal},$$

where $\bar{x} = \frac{1}{N}X\mathbf{1}$, $\bar{u} = \frac{1}{N}U\mathbf{1}$, $\bar{x}^+ = \frac{1}{N}X^+\mathbf{1}$ and $\bar{\xi} = \frac{1}{N}\sum_{k=1}^N \xi^k$.

Proposition

Consider a linear stochastic system and $v \in \mathbb{V}_q$. Then, given a dataset (X, U, X^+) ,

$$G\left(ar{x},ar{u},ar{\xi}
ight)=ar{x}ar{x}^{\intercal}-\gammaar{x}^{+}ar{x}^{+\intercal},$$

where $\bar{x} = \frac{1}{N}X1$, $\bar{u} = \frac{1}{N}U1$, $\bar{x}^+ = \frac{1}{N}X^+1$ and $\bar{\xi} = \frac{1}{N}\sum_{k=1}^N \xi^k$.

• G can potentially be computed for any (x, u) pair depending on the available data

Proposition

Consider a linear stochastic system and $v \in \mathbb{V}_q$. Then, given a dataset (X, U, X^+) ,

$$G\left(ar{x},ar{u},ar{\xi}
ight)=ar{x}ar{x}^{\intercal}-\gammaar{x}^{+}ar{x}^{+\intercal},$$

where $\bar{x} = \frac{1}{N}X\mathbf{1}$, $\bar{u} = \frac{1}{N}U\mathbf{1}$, $\bar{x}^+ = \frac{1}{N}X^+\mathbf{1}$ and $\bar{\xi} = \frac{1}{N}\sum_{k=1}^N \xi^k$.

G can potentially be computed for any (x, u) pair depending on the available data
 G(x, u, ξ) ≈ G(x, u, 0) = H(x, u)

Proposition

Consider a linear stochastic system and $v \in \mathbb{V}_q$. Then, given a dataset (X, U, X^+) ,

$$G\left(ar{x},ar{u},ar{\xi}
ight)=ar{x}ar{x}^{\intercal}-\gammaar{x}^{+}ar{x}^{+\intercal},$$

where $\bar{x} = \frac{1}{N}X\mathbf{1}$, $\bar{u} = \frac{1}{N}U\mathbf{1}$, $\bar{x}^+ = \frac{1}{N}X^+\mathbf{1}$ and $\bar{\xi} = \frac{1}{N}\sum_{k=1}^N \xi^k$.

- G can potentially be computed for any (x, u) pair depending on the available data
- $G(x, u, \overline{\xi}) \approx G(x, u, 0) = H(x, u)$
- ► The approximation is: $vec(H(x, u) \gamma Z)^{\intercal} vec(P) \leq \ell(x, u)$

Proposition

Consider a linear stochastic system and $v \in \mathbb{V}_q$. Then, given a dataset (X, U, X^+) ,

$$G\left(ar{x},ar{u},ar{\xi}
ight)=ar{x}ar{x}^{\intercal}-\gammaar{x}^{+}ar{x}^{+\intercal},$$

where $\bar{x} = \frac{1}{N}X\mathbf{1}$, $\bar{u} = \frac{1}{N}U\mathbf{1}$, $\bar{x}^+ = \frac{1}{N}X^+\mathbf{1}$ and $\bar{\xi} = \frac{1}{N}\sum_{k=1}^N \xi^k$.

• G can potentially be computed for any (x, u) pair depending on the available data

•
$$G(x, u, \overline{\xi}) \approx G(x, u, 0) = H(x, u)$$

- ► The approximation is: $vec(H(x, u) \gamma \Sigma)^{\intercal} vec(P) \leq \ell(x, u)$
- Optimal value function is shifted but the optimal policy is preserved!



Conclusions and future work

In summary:

- Under LQ assumptions, a sufficiently rich dataset can be used to artificially generate all Bellman inequalities
- ► The associated stage-cost can be reconstructed thanks to a bilinear algebra framework
- In case of stochastic systems, we can provide an (intentionally) biased estimate of the Bellman inequalities that preserves the optimal policy without iterative re-initialization

Conclusions and future work

In summary:

- Under LQ assumptions, a sufficiently rich dataset can be used to artificially generate all Bellman inequalities
- ► The associated stage-cost can be reconstructed thanks to a bilinear algebra framework
- In case of stochastic systems, we can provide an (intentionally) biased estimate of the Bellman inequalities that preserves the optimal policy without iterative re-initialization

Further developments:

- How to synthesise (approximate) Bellman inequalities from data for stochastic systems
- Relax the LQ assumptions to polynomial





