

# LQ Optimal Tracking with Unbounded Cost for Unknown Dynamics: An Adaptive Dynamic Programming Approach

Sebastian Bernhard and Jürgen Adamy

**Abstract**—In case of unknown system dynamics, we consider linear quadratic optimal tracking on infinite horizons with generally unbounded cost. For the first time, we deal with this problem in the framework of adaptive dynamic programming. So far, existing methods require bounded costs which essentially limits the applicability and achievable performance. Thus, we develop a new algorithm that yields a strongly overtaking optimal control which is an adequate solution. After collecting measurement data in an exploration phase, the algorithm implicitly solves the necessary and sufficient algebraic equations in [3], but without knowledge of the dynamics. Then, implementing the control results in an optimal transition to an optimal stationary trajectory. A simulation example of almost exact tracking for an over-actuated system demonstrates a highly efficient saving of input-energy in contrast to state-of-the-art approaches.

## I. INTRODUCTION

In the recent past, methods for online learning of optimal controls for linear systems with unknown dynamics have become a matter of great interest [15], [16]. Concentrating on continuous-time systems, first, the essential linear quadratic regulator (LQR) problem was solved in [12], [21] by a technique which is known as *adaptive dynamic programming* [15]. From there on, several approaches achieved exact tracking [8], [9] and infinite-horizon linear quadratic optimal tracking (LQT) [17], [18], [20]. So far, the latter only handle LQT problems (LQTP) for the very limited case of *bounded* cost. Before we present a literature overview and discuss the advantages of an LQTP that allows for *unbounded* cost, let us state our **main problem** and our key contributions.

In this paper, we regard a disturbed LTI system

$$\dot{x} = Ax + Bu + E_d \bar{x}, \quad (1a)$$

$$y = Cx + D_d \bar{x} \quad (1b)$$

with states  $x \in \mathbb{R}^n$ ,  $x(0) = x_0$ , inputs  $u \in \mathbb{R}^m$ , outputs  $y \in \mathbb{R}^p$  and state/output disturbances  $E_d \bar{x}$ ,  $D_d \bar{x}$  given by a *not stable* exosystem with states  $\bar{x} \in \mathbb{R}^{\bar{n}}$ ,  $\bar{x}(0) = \bar{x}_0$ :

$$\dot{\bar{x}} = \bar{A} \bar{x}, \quad (2a)$$

$$\bar{y} = \bar{C} \bar{x} \quad (2b)$$

which also generates the reference output  $\bar{y} \in \mathbb{R}^p$ . We pose

**Linear Quadratic Tracking Problem 1:** For systems (1) and (2), find the optimal control  $u^*(\cdot)$  concerning the cost

$$J_{t_E, T}(u(\cdot)) \Big|_{x(t_E), \bar{x}(t_E)} = \frac{1}{2} \int_{t_E}^T (y - \bar{y})^\top Q (y - \bar{y}) + u^\top R u \, dt, \quad (3)$$

The authors are with the Institute of Automatic Control and Mechatronics, Control Methods and Robotics Lab; Technische Universität Darmstadt, Landgraf-Georg Str. 4, 64283 Darmstadt, Germany, {bernhard, adamy}@rnr.tu-darmstadt.de

with  $t_E \geq 0$ , constant  $Q \succeq 0$  and  $R \succ 0$ , on an infinite horizon  $T \rightarrow \infty$  when the *data* of the complete dynamics:

$$A, B, E_d \text{ and } \bar{A} \quad (4)$$

is **unknown**.

Even for known data, this is in effect a non-trivial problem for two reasons. First, the cost is unbounded in general, i.e.  $\lim_{T \rightarrow \infty} J_{t_E, T}(u(\cdot)) \Big|_{x(t_E), \bar{x}(t_E)} = +\infty$  for any  $u(\cdot)$  (see [1]), which requires an adequate concept of optimality [6]. Second, there are no suitable sufficient conditions to prove optimality for the given LQTP, cf. [7].

To our advantage, the two mentioned problems have been rigorously solved [2], [3]. In contrast to [2], a time-invariant control of simple structure  $\hat{u} = -\widehat{K}(\hat{x} - \Pi_x^* \bar{x}) + F^* \bar{x}$  is derived in [3] based on necessary optimality criteria for  $T \rightarrow \infty$  [10]. Under certain conditions, it is *strongly overtaking optimal* with respect to cost (3), see [2] for a definition. As a result, the transition to the stationary state:  $\lim_{t \rightarrow \infty} \hat{x}(t) - \Pi_x^* \bar{x}(t) = 0$  is optimal with LQR feedback  $\widehat{K}$ . Moreover, it is proven in [3] that  $\Pi_x^* \bar{x}$  is an *optimal stationary state*. Thus, any other  $u(\cdot) \neq \hat{u}(\cdot)$  leading to  $x(t) - \Pi_x^* \bar{x}(t) \not\rightarrow 0$  as  $t \rightarrow \infty$  yields:  $\lim_{T \rightarrow \infty} J_{t_E, T}(u(\cdot)) - J_{t_E, T}(\hat{u}(\cdot)) = +\infty$  and, therefore, it requires infinitely more cost. Following the concept of *agreeable plans* [6],  $\hat{u}(\cdot)$  is also shown to be a favorable approximation of the optimal control  $u_T^*(\cdot)$  on finite horizons  $[t_E, T]$  under relaxed assumptions.

Considering the transition, it is beneficial to be able to specify a different cost than (3). It allows us to handle differing or even conflicting requirements of stationary behavior and transition, cf. [3]. Consequently, we have to implement

$$u^* = -K^*(x^* - \Pi_x^* \bar{x}) + F^* \bar{x} = -K^* \tilde{x}^* + F^* \bar{x} \quad (5)$$

where control  $\tilde{u}^* = -K^* \tilde{x}^*$  solves a separately formulated LQR problem (LQRP) whereas  $\lim_{t \rightarrow \infty} x^*(t) - \Pi_x^* \bar{x}(t) = 0$  is still guaranteed. However, to obtain  $K^*$ ,  $\Pi_x^*$  and  $F^*$ , it is required to solve necessary and sufficient algebraic equations in [3] which depend on the unknown data (4). In this light, we formulate our **main contributions**:

For the first time, we develop an adaptive dynamic programming (ADP) approach to determine control (5) in three steps: First, measurement data is collected during an exploration phase  $[0, t_E]$ . Second,  $u^*(\cdot)$  is obtained by calculating

C1) a feedback matrix  $K_k$ ,  $k \in \mathbb{N}_0$ , iteratively such that  $\lim_{k \rightarrow \infty} \|K_k - K^*\|_2 = 0$  as well as  $\Pi_x^*$  and  $F^*$  exactly such that  $\lim_{t \rightarrow \infty} x^*(t) - \Pi_x^* \bar{x}(t) = 0$  *without* requiring the unknown data (4).

C2) Third, implementing  $u^*(\cdot)$  on  $[t_E, \infty)$  guarantees an optimal transition to stationary state  $\Pi_x^* \bar{x}$ , which is

optimal on infinite horizons and approximates optimal solutions on finite horizons, under certain assumptions.

In addition,

- C3)  $\mathbf{K}^*$  may solve a separated LQR problem in order to consider a transition cost different from (3) and to obtain an optimal transition  $\|\tilde{\mathbf{x}}^*(t)\|_2 \leq M e^{-\alpha t}$ ,  $M \in \mathbb{R}^{>0}$  with a specified degree of stability  $\alpha > 0$ .
- C4) The approach is well suited in case of *under-actuation*  $m < p$ , where  $\lim_{t \rightarrow \infty} \mathbf{y}(t) - \bar{\mathbf{y}}(t) = \mathbf{0}$  is usually infeasible, and *over-actuation*  $\text{rank}(\mathbf{B}) > p$  without any need of solving additional parametric optimization problems in contrast to, e.g., [8], [9].

In the sequel, we give a short **literature review** and highlight some advantages of our approach.

In the domain of interest, the first results considered the online identification of the LQR feedback  $\mathbf{K}^*$  of a plant (1a),  $\bar{\mathbf{x}}(\cdot) \equiv \mathbf{0}$ , with unknown system matrix  $\mathbf{A}$  [21] and additionally unknown input matrix  $\mathbf{B}$  [12]. The proposed methods exploit the iterative computation of the *Riccati matrix* given by [13] and involve two separated steps. First, collecting measurement data in a phase where the system is excited by exploration noise and stabilized by an initial feedback  $\mathbf{K}_0$ . Second, computing iteratively a  $\mathbf{K}_k$  by a procedure equivalent to [13] but independent of unknown data (4) such that convergence occurs:  $\lim_{k \rightarrow \infty} \mathbf{K}_k = \mathbf{K}^*$ . This technique known as adaptive dynamic programming (ADP) or, at times, (measurement-) *data-driven adaptive optimal control* is used to derive the first part of C1). Moreover, we extend the results of [12] to account for C3), i.e. guaranteeing a specified degree of stability. This is useful for preserving the degree of stability given by  $\mathbf{K}_0$ .

Regarding the problem of tracking references, all results [8], [9] and [17], [18], [20] consider the same setup as ours, i.e. LTI-systems (1a) with references (and disturbances) generated by an exosystem (2a). We separate two groups.

Based on output regulation theory, [8] and [9] introduce exact tracking for unknown dynamics without steady-state error:  $\lim_{t \rightarrow \infty} \mathbf{y}(t) - \bar{\mathbf{y}}(t) = \mathbf{0}$ . This is achieved by an LQR-transition, obtained via ADP methods as above, to a stationary state that is the solution of the necessary and sufficient *regulator equations* [19]. Before solving these, the parts depending on unknown matrices have to be identified. In [8], [9], this requires to solve repeatedly a set of equations for matrix variables. In addition, if (1a) is over-actuated, i.e.  $\text{rank}(\mathbf{B}) > \text{rank}(\mathbf{C})$ , the sought solution is not unique. Then, an optimization problem by [14] is considered in [8], [9] which has to be solved online for  $n\bar{n} + m\bar{n}$  variables. In any case, we require to solve a single set of equations only once and, thus, our implementation is less complicated. Furthermore, we show by simulation that our approach is more efficient in using additional actuators.

The second group [17], [18], [20] considers LQT problems on infinite horizons which are essentially limited to bounded cost. All of them merge (1a), (2a) into an augmented system with state  $\mathbf{x}_{\text{aug}}^T = [\mathbf{x}^T \ \bar{\mathbf{x}}^T]$ . Then, they solve LQTP 1 as an LQRP of the augmented system. Indeed, this is a common

approach for the so-called servo-problem [1, Sec. 4.2]. In [18], LQTP 1 is solved by the methods in [12] in order to obtain the desired LQR  $\mathbf{u}_{\text{aug}}^*(\cdot)$ . This requires, however, that a  $\mathbf{P}_{\text{aug}} \in \mathbb{R}^{(n+\bar{n}) \times (n+\bar{n})}$  exists such that  $J_{t_E, \infty}(\mathbf{u}_{\text{aug}}^*(\cdot)) = \frac{1}{2} \mathbf{x}_{\text{aug}}(t_E)^T \mathbf{P}_{\text{aug}} \mathbf{x}_{\text{aug}}(t_E)$ , i.e. the cost is *bounded*. Following [1], this is only true when the augmented system is stabilizable. Necessarily, [18] assumes that  $\bar{\mathbf{A}}$  is *Hurwitz*, i.e.  $\bar{\mathbf{y}}(\cdot)$  must be exponentially decreasing which is a crucial restriction. The same assumption is obviously needed to apply the results of [20]. There, an actor-critic structure (cf. [16]) uses a critic to improve the actor. The critic precisely mimics the same quadratic form  $\frac{1}{2} \mathbf{x}_{\text{aug}}^T \mathbf{P}_{\text{aug}} \mathbf{x}_{\text{aug}}$  via a  $Q$ -function (cf. [20, Lemma 2]) which necessarily needs to exist. In contrast we assume, without loss of generality, that all eigenvalues of  $\bar{\mathbf{A}}$  lie in the closed right half-plane.

To handle the problem above, [17] uses a discount-factor  $\gamma > 0$  by multiplying the integrand of (3) by  $e^{-2\gamma t}$ . Similar to [1, Sec. 3.5], the discounted LQRP for the augmented system  $\mathbf{x}_{\text{aug}}$  can be reformulated as a standard LQRP. Then again, the solution in [17] can be obtained by using methods in [12], [21] if  $\gamma > 0$  is chosen such that  $\bar{\mathbf{A}} - \gamma \mathbf{I}$  is *Hurwitz*. That is a milder assumption as above. Using a discount factor, however, leads to several shortcomings which are absent in our approach. First, the discount factor lowers the degree of stability of the closed-loop dynamics in an opposite manner as in [1, Sec. 3.5]. One can observe that it may cause unstable closed loops when the open loop is unstable. Second, it is rather not-transparent how the choice of  $\gamma > 0$  affects the stationary tracking performance.

Altogether, [17], [18], [20] have to use the same cost for transition and stationary behavior. Then for “large” weights  $\mathbf{Q}$ , input constraints may be violated during transition. This is critical, if “large” weights are desired for *almost exact tracking* (AET), i.e. implementing arbitrarily small stationary tracking errors on any finite horizon of interest. This motivates C4) that allows us to implement AET more efficiently for over-actuated systems than [8], [9], cf. Section V.

Summarizing, there is strong motivation to solve LQTP 1 for generally unbounded cost in case of unknown dynamics, i.e. without restrictive assumptions [18] or modifications of (3) as [17]. To achieve C1-4), we proceed as follows. After some preliminaries, we solve our problem for (4) assumed known in Section III. In Section IV, we modify ADP in [9], [12] to account for C3). Then, we present a new set of equations which is solved for  $\Pi_x^*$ ,  $\mathbf{F}^*$  in view of C1-2). An online algorithm and a simulation demonstrating C4) are provided in Section V. Finally, we give concluding remarks.

## II. MATHEMATICAL PRELIMINARIES

We start by defining operations which are used in ADP, e.g. [8], [12], [21]. For a matrix  $\mathbf{P} = \mathbf{P}^T \in \mathbb{R}^{n \times n}$ , the bijective mapping  $\hat{\mathbf{P}} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{0.5n(n+1)}$  is defined as

$$\hat{\mathbf{P}} := [p_{1,1} \ 2p_{1,2} \ \cdots \ 2p_{1,n} \ p_{2,2} \ 2p_{2,3} \ \cdots \ 2p_{n-1,n} \ p_{n,n}]^T \quad (6)$$

where  $p_{i,j}$  is the element in the  $i$ -th row and  $j$ -th column of  $\mathbf{P}$ . For vector  $\mathbf{x} \in \mathbb{R}^n$ , we define  $\underline{\mathbf{x}} : \mathbb{R}^n \rightarrow \mathbb{R}^{0.5n(n+1)}$  as

$$\underline{\mathbf{x}} := [x_1^2 \ x_1 x_2 \ \cdots \ x_1 x_n \ x_2^2 \ x_2 x_3 \ \cdots \ x_{n-1} x_n \ x_n^2]^T \quad (7)$$

Furthermore, we will make use of the *Kronecker product* denoted by  $\mathbf{A} \otimes \mathbf{B} : \mathbb{R}^{m \times n} \times \mathbb{R}^{o \times p} \rightarrow \mathbb{R}^{mo \times np}$  for which the following relations hold:

$$(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top, \quad (8a)$$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}, \quad (8b)$$

$$\mathbf{L}^{(o,m)}(\mathbf{A} \otimes \mathbf{B})\mathbf{L}^{(n,p)} = \mathbf{B} \otimes \mathbf{A}, \quad (8c)$$

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A})\text{vec}(\mathbf{B}), \quad (8d)$$

$$\mathbf{x}^\top \mathbf{P} \mathbf{x} = \underline{\mathbf{x}}^\top \hat{\mathbf{P}} = (\mathbf{x}^\top \otimes \mathbf{x}^\top) \mathbf{T}^\top \hat{\mathbf{P}}, \quad (8e)$$

where we used  $\text{vec}(\mathbf{A}) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{nm}$  that stacks the columns of  $\mathbf{A}$  and the commutation matrix  $\mathbf{L}^{(q,r)} = \sum_{i=1}^q \sum_{j=1}^r (\mathbf{e}_{i,q} \mathbf{e}_{j,r}^\top) \otimes (\mathbf{e}_{j,r} \mathbf{e}_{i,q}^\top)$  with the unit vector  $\mathbf{e}_{k,l}$  of length  $l$  with  $k$ -th element equal one, zero else. Notice that in (8e) the mappings (6), (7) and  $\underline{\mathbf{x}} = \mathbf{T}(\mathbf{x} \otimes \mathbf{x})$  were utilized, where each row of  $\mathbf{T}$  is given by a specific  $\mathbf{e}_{i,n^2}^\top$  according to (7) and  $\mathbf{T}$  has full row rank.

We also use the *notations*: The zero and identity matrix of specified dimensions by  $\mathbf{0}_{a \times b}$  and  $\mathbf{I}_a$ , a positive (semi-) definite matrix  $\mathbf{M}$  by  $\mathbf{M} \succ (\succeq) \mathbf{0}$ , the spectrum of a matrix  $\mathbf{M}$  by  $\sigma(\mathbf{M})$ , 2-norm:  $\|\cdot\|_2$  and Frobenius norm:  $\|\cdot\|_F$ .

### III. TIME-INVARIANT CONTROL IN LINEAR QUADRATIC TRACKING

Based on [3], we derive control  $\mathbf{u}^*(\cdot)$  accounting for C1-3) if the data (4) is known. Besides, Kleinman's algorithm [13] to solve *algebraic Riccati equations* (ARE) is stated.

In view of C3), we want to consider different costs for transition and stationary behavior as suggested by [3]. Exploiting the structure of  $\mathbf{u}^*(\cdot)$  given by (5), we rewrite the corresponding solution of (1a) as  $\mathbf{x} = \tilde{\mathbf{x}} + \mathbf{\Pi}_x^* \bar{\mathbf{x}}$  dividing the stationary behavior  $\mathbf{\Pi}_x^* \bar{\mathbf{x}}$  that is determined by

$$\mathbf{\Pi}_x^* \bar{\mathbf{A}} = \mathbf{A} \mathbf{\Pi}_x^* + \mathbf{B} \mathbf{F}^* + \mathbf{E}_d$$

from the transition  $\tilde{\mathbf{x}}$  defined by the transient dynamics

$$\dot{\tilde{\mathbf{x}}} = \mathbf{A} \tilde{\mathbf{x}} + \mathbf{B} \tilde{\mathbf{u}} \quad (9)$$

with  $\tilde{\mathbf{x}}_0 = \mathbf{x}_0 - \mathbf{\Pi}_x^* \bar{\mathbf{x}}_0$  and  $\tilde{\mathbf{u}} = \mathbf{u} - \mathbf{F}^* \bar{\mathbf{x}}$ . This allows us to formulate an LQRP which accounts for an optimal transition:

**Linear Quadratic Regulator Problem 2:** For given  $\alpha \geq 0$  and dynamics (9), find the control  $\tilde{\mathbf{u}}^*(\cdot)$  which minimizes

$$\tilde{J}(\tilde{\mathbf{u}}(\cdot))|_{\tilde{\mathbf{x}}(t_E)} = \frac{1}{2} \int_{t_E}^{\infty} e^{2\alpha t} \left( \tilde{\mathbf{x}}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}} + \tilde{\mathbf{u}}^\top \tilde{\mathbf{R}} \tilde{\mathbf{u}} \right) dt, \quad (10)$$

where  $\tilde{\mathbf{Q}} \succeq \mathbf{0}$ ,  $\tilde{\mathbf{R}} \succ \mathbf{0}$  are constant and  $t_E \geq 0$ .

At all times, we consider the following assumptions

*Assumption 1:*  $(\mathbf{A} + \alpha \mathbf{I}, \mathbf{B})$  is stabilizable.

*Assumption 2:* For some  $\mathbf{V}^\top \mathbf{V} = \mathbf{C}^\top \mathbf{Q} \mathbf{C}$ ,  $(\mathbf{V}, \mathbf{A})$  is detectable and for some  $\mathbf{W}^\top \mathbf{W} = \tilde{\mathbf{Q}}$ ,  $(\mathbf{W}, \mathbf{A})$  is observable.

*Assumption 3:* It holds  $\sigma(\bar{\mathbf{A}}) \subset \{\lambda \in \mathbb{C} \mid \text{Re}(\lambda) \geq 0\}$ .

*Assumption 4:* The characteristic polynomial of  $\bar{\mathbf{A}}$  is a minimal polynomial (for a definition see, e.g., [19, p. 68]).

Asmp. 1 and 2 are standard in LQTP. The first is true  $\forall \alpha \geq 0$  if  $(\mathbf{A}, \mathbf{B})$  is controllable. The second implies that  $\forall \alpha \geq 0$ ,  $(\mathbf{W}, \mathbf{A} + \alpha \mathbf{I})$  is observable. In view of Asmp. 3, we simply include any asymptotically stable part of (2a) into

(1a). Asmp. 4 specifies that the geometric multiplicity of each distinct eigenvalue of  $\bar{\mathbf{A}}$  equals one. It is also needed in, e.g., [8], [9], [17], [18], and we will remark where it is required.

We give a result extending [3, Coroll. 10], where we firstly consider *bounded* references in view of stationary tracking:

*Lemma 1:* For any given  $\alpha \geq 0$  and if  $\forall \bar{\mathbf{x}}_0 \exists M \in \mathbb{R}^{>0}$  such that  $\|\bar{\mathbf{x}}(t)\|_2 \leq M, \forall t \in [0, \infty)$ , then for the control

$$\mathbf{u}^* = \underbrace{-\mathbf{R}^{-1} \mathbf{B}^\top \mathbf{P}(\mathbf{x}^* - \mathbf{\Pi}_x^* \bar{\mathbf{x}})}_{-\mathbf{K}^* \tilde{\mathbf{x}}^*} \underbrace{- \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{\Pi}_\phi^* \bar{\mathbf{x}}}_{+\mathbf{F}^* \bar{\mathbf{x}}}, \quad (11)$$

where *Riccati* matrix  $\mathbf{P} = \mathbf{P}^\top \succ \mathbf{0}$  uniquely solves ARE

$$\mathbf{P}(\mathbf{A} + \alpha \mathbf{I}) + (\mathbf{A} + \alpha \mathbf{I})^\top \mathbf{P} - \mathbf{P} \mathbf{B} \tilde{\mathbf{R}}^{-1} \mathbf{B}^\top \mathbf{P} = -\tilde{\mathbf{Q}} \quad (12)$$

and  $\mathbf{\Pi}_x^*, \mathbf{\Pi}_\phi^*$  with  $\mathbf{S} = \mathbf{C}^\top \mathbf{Q}(\bar{\mathbf{C}} - \mathbf{D}_d)$  uniquely solve

$$\begin{bmatrix} \mathbf{\Pi}_x \\ \mathbf{\Pi}_\phi \end{bmatrix} \bar{\mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{A} & -\mathbf{B} \tilde{\mathbf{R}}^{-1} \mathbf{B}^\top \\ -\mathbf{C}^\top \mathbf{Q} \mathbf{C} & -\mathbf{A}^\top \end{bmatrix}}_{=\mathbf{H}} \begin{bmatrix} \mathbf{\Pi}_x \\ \mathbf{\Pi}_\phi \end{bmatrix} + \begin{bmatrix} \mathbf{E}_d \\ \mathbf{S} \end{bmatrix}, \quad (13)$$

it is satisfied for any given  $\mathbf{x}(t_E) \in \mathbb{R}^n$  and  $\bar{\mathbf{x}}(t_E) \in \mathbb{R}^{\bar{n}}$ :

- 1) The pair  $(\tilde{\mathbf{x}}^*, \tilde{\mathbf{u}}^*)$  with  $\tilde{\mathbf{u}}^* = -\mathbf{K}^* \tilde{\mathbf{x}}^*$  is the optimal solution of LQRP 2. That is,  $\mathbf{u}^*$  leads to an **optimal transition**  $\lim_{t \rightarrow \infty} \tilde{\mathbf{x}}^* = \mathbf{0}$  with respect to cost (10).
- 2) The closed-loop matrix  $\mathbf{A} - \mathbf{B} \mathbf{K}^*$  is *Hurwitz* and the degree of stability  $\alpha$  is guaranteed, i.e.  $\max_{\lambda \in \sigma(\mathbf{A} - \mathbf{B} \mathbf{K}^*)} \text{Re}(\lambda) < -\alpha$ .
- 3) For any admissible control  $\mathbf{u}(\cdot) \neq \mathbf{u}^*(\cdot)$  on  $[t_E, \infty)$  such that  $\mathbf{x}(t) - \mathbf{\Pi}_x^* \bar{\mathbf{x}}(t) \not\rightarrow \mathbf{0}$  for  $t \rightarrow \infty$  it holds

$$\lim_{T \rightarrow \infty} J_{t_E, T}(\mathbf{u}(\cdot)) - J_{t_E, T}(\mathbf{u}^*(\cdot)) = +\infty.$$

That is,  $\mathbf{u}^*(\cdot)$  leads to an **optimal stationary solution** in the sense that any deviation from  $\mathbf{\Pi}_x^* \bar{\mathbf{x}}$  results in infinite additional cost.

*Proof:* 1) and 2) are standard results, c.f. [1, Sec. 3.5]. In view of 3), we extend [3, Coroll. 10]. Note that (13) parametrizes the particular solution of the *Hamiltonian* system which emerges from the necessary optimality conditions for LQTP 1 by [10]. Due to the uniqueness of solutions of (13), e.g. see [4, Thm. 3],  $\mathbf{u}^*(\cdot)$  and the unique optimal solution  $\hat{\mathbf{u}}(\cdot)$  of LQTP 1 derived in [3, Thm. 9] share the same stationary state  $\mathbf{\Pi}_x^* \bar{\mathbf{x}}$  and pre-filter  $\mathbf{F}^* \bar{\mathbf{x}}$ , i.e.  $\hat{\mathbf{u}}(\cdot) = -\tilde{\mathbf{K}}(\hat{\mathbf{x}} - \mathbf{\Pi}_x^* \bar{\mathbf{x}}) + \mathbf{F}^* \bar{\mathbf{x}}$ . Actually, if  $\tilde{\mathbf{Q}} = \mathbf{C}^\top \mathbf{Q} \mathbf{C}$ ,  $\tilde{\mathbf{R}} = \mathbf{R}$  and  $\alpha = 0$ , then  $\mathbf{u}^*(\cdot) \equiv \hat{\mathbf{u}}(\cdot)$ .

As a consequence, we have  $\lim_{t \rightarrow \infty} \mathbf{x}^*(t) - \hat{\mathbf{x}}(t) = \mathbf{0}$  and  $\lim_{t \rightarrow \infty} \mathbf{u}^*(t) - \hat{\mathbf{u}}(t) = \mathbf{0}$  both exponentially in  $t$ . By means of the *calculus of variations*, omitting the technicalities, it is rather evident to conclude that  $\lim_{T \rightarrow \infty} J_{t_E, T}(\mathbf{u}^*(\cdot)) - J_{t_E, T}(\hat{\mathbf{u}}(\cdot)) \leq M(\mathbf{x}(t_E), \bar{\mathbf{x}}(t_E))$ ,  $M(\mathbf{x}(t_E), \bar{\mathbf{x}}(t_E)) > 0$ . For showing this, one could study the first case in the given case study of [3, proof of Thm. 9]. Furthermore, we recall that it indeed holds  $\lim_{T \rightarrow \infty} J_{t_E, T}(\mathbf{u}(\cdot)) - J_{t_E, T}(\hat{\mathbf{u}}(\cdot)) = +\infty$ , see [3, Coroll. 10].

Using both, 3) readily follows after minor calculations which we leave to the reader. ■

In case of *unbounded* references and disturbances, 3) cannot be shown. But then, one is interested in a finite horizon rather

than  $T \rightarrow \infty$ . To this end, we adopt a result proven in [3, Thm. 12] which focuses on the stationary behavior here:

*Lemma 2:* Consider  $T_1 \geq 0$ :  $\mathbf{x}(T_1) = \mathbf{\Pi}_x^* \bar{\mathbf{x}}(T_1)$ . If

$$\operatorname{Re}(\lambda_H) + \operatorname{Re}(\bar{\lambda}) < 0 \quad (14)$$

for all  $\lambda_H \in \sigma(\mathbf{H}) \cap \{\lambda \in \mathbb{C} \mid \operatorname{Re}(\lambda) < 0\}$  and  $\bar{\lambda} \in \sigma(\bar{\mathbf{A}})$ , the pair  $(\mathbf{\Pi}_x^* \bar{\mathbf{x}}, \mathbf{F}^* \bar{\mathbf{x}})$  which solves (1a) on  $[T_1, \infty)$  satisfies

$$J_{T_1, \theta}(\mathbf{F}^* \bar{\mathbf{x}})|_{\bar{\mathbf{x}}(T_1)} - \lim_{T \rightarrow \infty} J_{T_1, \theta}(\mathbf{u}_T^*(\cdot))|_{\bar{\mathbf{x}}(T_1)} = 0 \quad (15)$$

for any  $\theta \geq T_1$  and any given  $\bar{\mathbf{x}}(T_1) \in \mathbb{R}^{\bar{n}}$ , where  $\mathbf{u}_T^*(\cdot)$  is the optimal control of LQTP 1 on a finite horizon  $[T_1, T]$ . ■ Clearly, (15) is a very desirable property which we put into words. Assume the transition phase ends after some  $T_1$ , i.e.  $\mathbf{x}^*(T_1) \approx \mathbf{\Pi}_x^* \bar{\mathbf{x}}(T_1)$ , and  $T$  is large, then the stationary behavior  $\mathbf{\Pi}_x^* \bar{\mathbf{x}}$  induced by  $\mathbf{u}^*(\cdot)$  closely approximates the optimal solution  $\mathbf{x}_T^*(\cdot)$  induced by  $\mathbf{u}_T^*(\cdot)$  on the interval  $[T_1, \theta]$ . That is, both produce approximatively the same cost as implied by (15). In this way,  $\mathbf{u}^*(\cdot)$  qualifies as an *agreeable plan* on  $[T_1, T]$ , cf. [6]. Meaning that for a large horizon  $T$ , which may not be exactly known, it is agreeable to apply  $\mathbf{u}^*(\cdot)$  instead of  $\mathbf{u}_T^*(\cdot)$  over a large interval.

*Remark 1:* Condition (14) needs only to be considered if  $\exists \bar{\lambda} \in \sigma(\bar{\mathbf{A}})$ :  $\operatorname{Re}(\bar{\lambda}) > 0$  and should be satisfied for a “large”  $\mathbf{Q}$  in most cases.

Our main focus lies on deriving (11) without the need of the unknown data (4). The associated cost of  $\mathbf{u}^*(\cdot)$  will serve as a starting point later on. It is concisely written as

$$J_{t, T}(\mathbf{u}^*(\cdot)) = \frac{1}{2} \mathbf{x}(t)^\top \mathbf{P} \mathbf{x}(t) + \mathbf{x}(t)^\top (-\mathbf{P} \mathbf{\Pi}_x^* + \mathbf{\Pi}_\phi^*) \bar{\mathbf{x}}(t) + z(\bar{\mathbf{x}}(t), \mathbf{x}(T), \bar{\mathbf{x}}(T), t, T) \quad (16)$$

if  $\tilde{\mathbf{Q}} = \mathbf{C}^\top \mathbf{Q} \mathbf{C}$ ,  $\tilde{\mathbf{R}} = \mathbf{R}$ ,  $\alpha = 0$ , some terms substituted by  $z$ . Note that (16) is generally unbounded as  $\lim_{T \rightarrow \infty} z = \infty$ .

Finally, let us introduce Kleinman’s algorithm [13] which lays the foundation for adaptive dynamic programming in, e.g., [9], [12], [21]. We adopt the results of [13] by minor modifications to our case of a specified degree of stability:

*Lemma 3:* For a given  $\alpha \geq 0$ , suppose  $\mathbf{K}_0$  such that  $(\mathbf{A} + \alpha \mathbf{I} - \mathbf{B} \mathbf{K}_0)$  is *Hurwitz*. Consider

$$\mathbf{P}_k \mathbf{A}_k + \mathbf{A}_k^\top \mathbf{P}_k + 2\alpha \mathbf{P}_k = -\tilde{\mathbf{Q}}_k, \quad (17a)$$

$$\text{For } k \geq 1: \mathbf{K}_k = \tilde{\mathbf{R}}^{-1} \mathbf{B}^\top \mathbf{P}_{k-1}, \quad (17b)$$

with  $\mathbf{A}_k = \mathbf{A} - \mathbf{B} \mathbf{K}_k$  and  $\tilde{\mathbf{Q}}_k = \tilde{\mathbf{Q}} + \mathbf{K}_k^\top \tilde{\mathbf{R}} \mathbf{K}_k$  then for any  $k \in \mathbb{N}_0$  the following holds:

- 1)  $(\mathbf{A} + \alpha \mathbf{I} - \mathbf{B} \mathbf{K}_k)$  is *Hurwitz*. As a consequence,  $\max_{\lambda \in \sigma(\mathbf{A} - \mathbf{B} \mathbf{K}_k)} \operatorname{Re}(\lambda) < -\alpha$ .
- 2)  $\mathbf{P}_k$  is symmetric and  $\mathbf{P}_k \succeq \mathbf{P}_{k+1} \succeq \mathbf{P} \succ \mathbf{0}$  where  $\mathbf{P}$  uniquely solves (12).
- 3)  $\lim_{k \rightarrow \infty} \mathbf{P}_k = \mathbf{P}$ . ■

#### IV. APPLICATION TO UNKNOWN DYNAMICS: AN ADP APPROACH

For the first time, we achieve C2) despite the unknown data (4). That is, implementing  $\mathbf{u}^*(\cdot)$  as in (11) with desired properties as stated in Lemma 1 and 2. Apparently, we have to find  $\mathbf{K}^*$ ,  $\mathbf{\Pi}_x^*$  and  $\mathbf{F}^*$ .

To this end, we present C1). Firstly, we show how to obtain  $\mathbf{K}_k \approx \mathbf{K}^*$  iteratively. In the process, unknown  $\mathbf{B}$  and  $\mathbf{E}_d$  are also provided. It is based on adaptive dynamic programming (ADP) [9], [12]. We extend these methods to achieve a desired degree of stability C3). Secondly, we introduce a new approach which allows us to determine  $\mathbf{\Pi}_x^*$  and  $\mathbf{F}^*$  exactly in spite of the unknown  $\mathbf{A}$ ,  $\bar{\mathbf{A}}$  and approximation  $\mathbf{P}_k \approx \mathbf{P}$ .

#### A. Optimal Transition of Unknown Dynamics

Our goal is to solve LQRP 2 despite the unknown data (4). LQRPs for  $\alpha = 0$  were approached in [9], [12] and [21] by the method of ADP. Following [12], the idea is to collect *measurement data* containing sufficient information on the unknown dynamics (4). Thus, we collect measurement data depending on  $\mathbf{x}(t)$ ,  $\mathbf{u}(t)$  and  $\bar{\mathbf{x}}(t)$  during an exploration phase  $[0, t_E]$  where the system (1a) is excited by

$$\mathbf{u} = -\mathbf{K}_0 \mathbf{x} + \mathbf{e} \quad (18)$$

with  $\mathbf{e}$  being a *known* exploration noise. Afterwards, a set of equations must be provided which is equivalent to (17) but depends only on the collected measurement data instead of unknown structural data (4). Then, iteratively solving these equations gives  $\mathbf{K}_k \rightarrow \mathbf{K}^*$  for  $k \rightarrow \infty$  as desired.

*Remark 2:* When the open loop is not stable, requiring  $\mathbf{A} - \mathbf{B} \mathbf{K}_0$  being *Hurwitz* is a crucial step in Lemma 3 and, thus, in [9], [12], etc.. As suggested by [8], some information on the system parameters is usually available. Suppose we know convex polytopes  $\mathcal{A}$  and  $\mathcal{B}$ , then an optimization-based design [5] can be used to obtain a stabilizing  $\mathbf{K}_0$  for any  $\mathbf{A} \in \mathcal{A}$ ,  $\mathbf{B} \in \mathcal{B}$  and to quantify a minimal degree of stability  $\alpha_0$ . This motivates our extension of the results in [9], [12], i.e. we can specify  $0 \leq \alpha < \alpha_0$  in LQRP 2. By choosing  $\alpha$  close to  $\alpha_0$ , advantageously, any  $\mathbf{K}_k$  will preserve the degree of stability  $\alpha \approx \alpha_0$  independent of our choice of  $\tilde{\mathbf{Q}}$ ,  $\tilde{\mathbf{R}}$ .

Subsequently, we derive the desired set of equations and give a unique solution. We follow the derivations in [9], [12] and extend these for the case  $\alpha \geq 0$ . Let us denote  $\mathbf{x}(t)$  given by (1a) for control law (18). We integrate the derivative of  $\mathbf{x}(t)^\top \mathbf{P}_k \mathbf{x}(t)$ ,  $k \in \mathbb{N}_0$ , over some period  $[t_i, t_{i+1}] \subset [0, t_E]$  with  $t_{i+1} \geq t_i$ ,  $i \in \mathbb{N}$  and  $t_1 = 0$ . It holds

$$\begin{aligned} & \mathbf{x}^\top(t_{i+1}) \mathbf{P}_k \mathbf{x}(t_{i+1}) - \mathbf{x}^\top(t_i) \mathbf{P}_k \mathbf{x}(t_i) + 2\alpha \int_{t_i}^{t_{i+1}} \mathbf{x}^\top \mathbf{P}_k \mathbf{x} \, d\tau \\ &= \int_{t_i}^{t_{i+1}} \mathbf{x}^\top \left( \mathbf{A}_k^\top \mathbf{P}_k + \mathbf{P}_k \mathbf{A}_k + 2\alpha \mathbf{P}_k + 2(\mathbf{K}_k \mathbf{x})^\top \mathbf{B}^\top \mathbf{P}_k \right) \mathbf{x} \\ & \quad + 2\mathbf{u}^\top \mathbf{B}^\top \mathbf{P}_k \mathbf{x} + 2(\mathbf{E}_d \bar{\mathbf{x}})^\top \mathbf{P}_k \mathbf{x} \, d\tau \\ &= \int_{t_i}^{t_{i+1}} -\mathbf{x}^\top \tilde{\mathbf{Q}}_k \mathbf{x} + 2(\mathbf{K}_k \mathbf{x} + \mathbf{u})^\top \tilde{\mathbf{R}} \mathbf{K}_{k+1} \mathbf{x} + 2\bar{\mathbf{x}}^\top \mathbf{E}_d^\top \mathbf{P}_k \mathbf{x} \, d\tau \end{aligned} \quad (19)$$

where we used (17a) and  $\mathbf{B}^\top \mathbf{P}_k = \tilde{\mathbf{R}} \mathbf{K}_{k+1}$  based on (17b).

At step  $k$ , (19) obviously depends on three unknown variables:  $\mathbf{K}_{k+1}$ , symmetric  $\mathbf{P}_k$  and the product  $\mathbf{E}_d^\top \mathbf{P}_k$ . Consequently, we follow the suggestion of [12] and regard  $i = 1, \dots, h$  time intervals such that  $[0, t_E] = \cup_{i=1}^h [t_i, t_{i+1}]$ . Then, stacking the  $h$  equations (19) leads to a set of equations which can be solved for the unknowns similarly as in [9].

With this in mind, (19) is reformulated by means of the mappings (6), (7) and relations (8), see Section II:

$$\begin{aligned} & \left( \mathbf{x}(t_{i+1}) - \mathbf{x}(t_i) + 2\alpha T \int_{t_i}^{t_{i+1}} (\mathbf{x} \otimes \mathbf{x}) \, d\tau \right)^\top \widehat{\mathbf{P}}_k \\ &= \int_{t_i}^{t_{i+1}} -(\mathbf{x}^\top \otimes \mathbf{x}^\top) \text{vec}(\widetilde{\mathbf{Q}}_k) + 2(\mathbf{x}^\top \otimes \bar{\mathbf{x}}^\top) \text{vec}(\mathbf{E}_d^\top \mathbf{P}_k) \\ & \quad + 2 \left( (\mathbf{x}^\top \otimes \mathbf{x}^\top) (\mathbf{I}_n \otimes \mathbf{K}_k^\top \widetilde{\mathbf{R}}) + (\mathbf{x}^\top \otimes \mathbf{u}^\top) (\mathbf{I}_n \otimes \widetilde{\mathbf{R}}) \right) \\ & \quad \cdot \text{vec}(\mathbf{K}_{k+1}) \, d\tau. \quad (20) \end{aligned}$$

Before we stack  $h$  equations (20), let us introduce

$$\Delta_{xx} = \begin{bmatrix} \mathbf{x}(t_2) - \mathbf{x}(t_1) & \cdots & \mathbf{x}(t_{h+1}) - \mathbf{x}(t_h) \end{bmatrix}^\top, \quad (21a)$$

$$\Lambda_{xx} = \begin{bmatrix} \int_{t_1}^{t_2} \mathbf{x} \otimes \mathbf{x} \, d\tau & \cdots & \int_{t_h}^{t_{h+1}} \mathbf{x} \otimes \mathbf{x} \, d\tau \end{bmatrix}^\top, \quad (21b)$$

$$\Lambda_{xu} = \begin{bmatrix} \int_{t_1}^{t_2} \mathbf{x} \otimes \mathbf{u} \, d\tau & \cdots & \int_{t_h}^{t_{h+1}} \mathbf{x} \otimes \mathbf{u} \, d\tau \end{bmatrix}^\top, \quad (21c)$$

$$\Lambda_{x\bar{x}} = \begin{bmatrix} \int_{t_1}^{t_2} \mathbf{x} \otimes \bar{\mathbf{x}} \, d\tau & \cdots & \int_{t_h}^{t_{h+1}} \mathbf{x} \otimes \bar{\mathbf{x}} \, d\tau \end{bmatrix}^\top \quad (21d)$$

with  $\Delta_{xx} \in \mathbb{R}^{h \times 0.5n(n+1)}$ ,  $\Lambda_{xx} \in \mathbb{R}^{h \times n^2}$ ,  $\Lambda_{xu} \in \mathbb{R}^{h \times nm}$ ,  $\Lambda_{x\bar{x}} \in \mathbb{R}^{h \times n\bar{n}}$ . Finally, one can obtain the set of equations

$$\Theta_k \begin{bmatrix} \widehat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{k+1}) \\ \text{vec}(\mathbf{E}_d^\top \mathbf{P}_k) \end{bmatrix} = -\Lambda_{xx} \text{vec}(\widetilde{\mathbf{Q}}_k) \quad (22)$$

with  $\Theta_k \in \mathbb{R}^{h \times (0.5n(n+1) + nm + n\bar{n})}$  given by

$$\Theta_k = \begin{bmatrix} \widetilde{\Delta}_{xx} & -2(\Lambda_{xx}(\mathbf{I}_n \otimes \mathbf{K}_k^\top \widetilde{\mathbf{R}}) + \Lambda_{xu}(\mathbf{I}_n \otimes \widetilde{\mathbf{R}})) & -2\Lambda_{x\bar{x}} \end{bmatrix}$$

where  $\widetilde{\Delta}_{xx} = \Delta_{xx} + 2\alpha \Lambda_{xx} T^\top$ . When  $\Theta_k$  has full column rank, clearly

$$\begin{bmatrix} \widehat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{k+1}) \\ \text{vec}(\mathbf{E}_d^\top \mathbf{P}_k) \end{bmatrix} = -(\Theta_k^\top \Theta_k)^{-1} \Theta_k^\top \Lambda_{xx} \text{vec}(\mathbf{Q}_k) \quad (23)$$

is the least square solution of (22). Existence can be guaranteed by a condition stated in the following lemma which is similar to [12, Lem. 6]:

*Lemma 4:* Suppose it holds

$$\text{rank} \left( \begin{bmatrix} \Lambda_{xx} & \Lambda_{xu} & \Lambda_{x\bar{x}} \end{bmatrix} \right) = \frac{1}{2}n(n+1) + nm + n\bar{n} \quad (24)$$

for the given  $[0, t_E] = \cup_{i=1}^h [t_i, t_{i+1}]$ , then, for any  $k \in \mathbb{N}_0$ ,  $\Theta_k$  has full column rank.

*Proof:*... goes along the same lines as [9, proof of Lem. 3 for  $i = 0$ ] with some modifications to cover  $\alpha \geq 0$ . ■

*Remark 3:* Obviously, the exploration noise  $e$  needs to be *sufficiently rich* in the sense that (24) holds. A typical choice is random noise or sums of sinusoidal signals [12]. From a practical perspective, the latter is preferable as it does not violate actuator rate constraints. Due to the structure of (21), we may need less frequencies than the suggested  $\frac{1}{2}(\frac{1}{2}n(n+1) + nm + n\bar{n})$  by [11, Sec. 5.2.1], cf. Section V.

Based on the derivation of (19), we know that, for any  $k \in \mathbb{N}_0$ , the solution  $\mathbf{P}_k, \mathbf{K}_{k+1}$  of (17) also satisfies (22). Hence, an exact solution of (22) exists which, of course, is given by the least square solution (23) when (24) holds.

As a consequence, iteratively solving (17) is equivalent to iteratively calculating (23). Following this discussion, the next theorem is self-evident. It extends [9, Thm. 2].

*Theorem 1:* For a given  $\alpha \geq 0$ , suppose  $\mathbf{K}_0$  such that  $\mathbf{A} + \alpha \mathbf{I} - \mathbf{B}\mathbf{K}_0$  is Hurwitz and condition (24) holds. For any  $k \in \mathbb{N}_0$ ,  $\mathbf{P}_k \succ \mathbf{0}$  and  $\mathbf{K}_{k+1}$  given by the policy iteration (23) satisfy (17) and it holds  $\lim_{k \rightarrow \infty} \mathbf{P}_k = \mathbf{P}$ ,  $\lim_{k \rightarrow \infty} \mathbf{K}_k = \mathbf{K}^*$  where  $\mathbf{P}, \mathbf{K}^*$  as in Lemma 1. ■

With  $\mathbf{P}_k$  invertible for any  $k \in \mathbb{N}_0$ , we can *exactly* calculate  $\mathbf{B} = \mathbf{P}_k^{-1} \mathbf{K}_{k+1}^\top \widetilde{\mathbf{R}}$  and  $\mathbf{E}_d$  from the third argument of (23) for some arbitrary  $k \in \mathbb{N}_0$ .

### B. Optimal Stationary Behavior of Unknown Dynamics

In the previous section, we obtained  $\mathbf{P}_k, \mathbf{K}_{k+1}$  as approximations of  $\mathbf{P}, \mathbf{K}^*$  and  $\mathbf{B}, \mathbf{E}_d$  exactly. In order to implement  $\mathbf{u}^*(\cdot)$ , it is left to determine  $\Pi_x^*$  and  $\mathbf{F}^* = -\mathbf{R}^{-1} \mathbf{B}^\top \Pi_\phi^*$ , i.e. we need the *exact* solution of (13) despite  $\mathbf{A}$  and  $\bar{\mathbf{A}}$  unknown. We present a new methodology. As a result, we solve LQTP 1 with unbounded cost in a stationarily optimal sense despite the unknown data for the first time.

As we see, for obtaining (19) a term similar to the quadratic, first addend in (16) was analyzed. In the sequel, we follow the idea of providing an analogous analysis for a term similar to the bilinear, second addend in (16).

First, we take  $\mathbf{P}_k$  for some  $k \in \mathbb{N}_0$  and consider

$$\begin{aligned} & \frac{d}{dt} \mathbf{x}^\top(t) \mathbf{P}_k \Pi_x \bar{\mathbf{x}}(t) \\ &= \underbrace{\mathbf{x}^\top (\mathbf{P}_k \Pi_x \bar{\mathbf{A}} + \mathbf{A}^\top \mathbf{P}_k \Pi_x) \bar{\mathbf{x}}}_{\mathbf{x}^\top \mathbf{P}_k (\Pi_x \bar{\mathbf{A}} - \mathbf{A} \Pi_x) \bar{\mathbf{x}} + \mathbf{x}^\top \Omega_k \Pi_x \bar{\mathbf{x}}} + (\mathbf{B}\mathbf{u} + \mathbf{E}_d \bar{\mathbf{x}})^\top \mathbf{P}_k \Pi_x \bar{\mathbf{x}} \quad (25) \end{aligned}$$

where we used (17a) to replace  $\mathbf{A}^\top \mathbf{P}_k = -\mathbf{P}_k \mathbf{A} + \Omega_k$  with

$$\Omega_k = -\widetilde{\mathbf{Q}}_k - 2\alpha \mathbf{P}_k + \mathbf{P}_k \mathbf{B} \mathbf{K}_k + \mathbf{K}_k^\top \mathbf{B}^\top \mathbf{P}_k.$$

After replacing  $\Pi_x \bar{\mathbf{A}} - \mathbf{A} \Pi_x$  in (25) with the help of the first line of (13), integration of (25) over an interval  $[t_i, t_{i+1}]$  as defined in Section IV-A yields

$$\begin{aligned} & \mathbf{x}^\top(t_{i+1}) \mathbf{P}_k \Pi_x \bar{\mathbf{x}}(t_{i+1}) - \mathbf{x}^\top(t_i) \mathbf{P}_k \Pi_x \bar{\mathbf{x}}(t_i) \\ &= \int_{t_i}^{t_{i+1}} -\mathbf{x}^\top \mathbf{P}_k \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \Pi_\phi \bar{\mathbf{x}} + \mathbf{x}^\top \mathbf{P}_k \mathbf{E}_d \bar{\mathbf{x}} \\ & \quad + \mathbf{x}^\top \Omega_k \Pi_x \bar{\mathbf{x}} + (\mathbf{B}\mathbf{u} + \mathbf{E}_d \bar{\mathbf{x}})^\top \mathbf{P}_k \Pi_x \bar{\mathbf{x}} \, d\tau. \quad (26) \end{aligned}$$

Second, we integrate the derivative of  $\mathbf{x}(t) \Pi_\phi \bar{\mathbf{x}}(t)$ :

$$\begin{aligned} & \mathbf{x}^\top(t_{i+1}) \Pi_\phi \bar{\mathbf{x}}(t_{i+1}) - \mathbf{x}^\top(t_i) \Pi_\phi \bar{\mathbf{x}}(t_i) \\ &= \int_{t_i}^{t_{i+1}} \mathbf{x}^\top (\Pi_\phi \bar{\mathbf{A}} + \mathbf{A}^\top \Pi_\phi) \bar{\mathbf{x}} + (\mathbf{B}\mathbf{u} + \mathbf{E}_d \bar{\mathbf{x}})^\top \Pi_\phi \bar{\mathbf{x}} \, d\tau \\ &= \int_{t_i}^{t_{i+1}} -\mathbf{x}^\top \mathbf{C}^\top \mathbf{Q} \mathbf{C} \Pi_x \bar{\mathbf{x}} + \mathbf{x}^\top \mathbf{S} \bar{\mathbf{x}} + (\mathbf{B}\mathbf{u} + \mathbf{E}_d \bar{\mathbf{x}})^\top \Pi_\phi \bar{\mathbf{x}} \, d\tau \quad (27) \end{aligned}$$

where we used the second line of (13).

Summarizing, we have found two equations (26) and (27) which are satisfied by the unique solution of (13) but do not depend on the unknown  $\mathbf{A}, \bar{\mathbf{A}}$ . Again as in Section IV-A, the idea is to obtain a set of equations equivalent to (13) by stacking (26), (27) for  $h$  time intervals. For this purpose, we

$$\bar{\Theta} = \begin{bmatrix} \Delta_{\bar{x}x}(\mathbf{I}_{\bar{n}} \otimes \mathbf{P}_k) - \Lambda_{\bar{x}x}(\mathbf{I}_{\bar{n}} \otimes \Omega_k) - \Lambda_{\bar{x}u}(\mathbf{I}_{\bar{n}} \otimes \mathbf{B}^\top \mathbf{P}_k) - \Lambda_{\bar{x}\bar{x}}(\mathbf{I}_{\bar{n}} \otimes \mathbf{E}_d^\top \mathbf{P}_k) & \Lambda_{\bar{x}x}(\mathbf{I}_{\bar{n}} \otimes \mathbf{P}_k \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top) \\ \Lambda_{\bar{x}x}(\mathbf{I}_{\bar{n}} \otimes \mathbf{C}^\top \mathbf{Q} \mathbf{C}) & \Delta_{\bar{x}x} - \Lambda_{\bar{x}u}(\mathbf{I}_{\bar{n}} \otimes \mathbf{B}^\top) - \Lambda_{\bar{x}\bar{x}}(\mathbf{I}_{\bar{n}} \otimes \mathbf{E}_d^\top) \end{bmatrix} \quad (29)$$

could rewrite (26), (27) in a similar manner as done in (20) by means of the relations (8). Due to the space restrictions, however, we omit these technical steps here and proceed directly to formulate the desired set of equations.

By stacking measurements for  $i = 1, \dots, h$ , we have

$$\Delta_{\bar{x}x} = [\dots \bar{x}(t_{i+1}) \otimes \mathbf{x}(t_{i+1}) - \bar{x}(t_i) \otimes \mathbf{x}(t_i) \dots]^\top, \quad (28a)$$

$$\Lambda_{\bar{x}x} = \Lambda_{\bar{x}x} \mathbf{L}^{(n, \bar{n})}, \quad (28b)$$

$$\Lambda_{\bar{x}u} = \left[ \int_{t_1}^{t_2} \bar{x} \otimes \mathbf{u} \, d\tau \quad \dots \quad \int_{t_h}^{t_{h+1}} \bar{x} \otimes \mathbf{u} \, d\tau \right]^\top, \quad (28c)$$

$$\Lambda_{\bar{x}\bar{x}} = \left[ \int_{t_1}^{t_2} \bar{x} \otimes \bar{x} \, d\tau \quad \dots \quad \int_{t_h}^{t_{h+1}} \bar{x} \otimes \bar{x} \, d\tau \right]^\top \quad (28d)$$

where we used  $\mathbf{L}^{(1,1)} = 1$  and with  $\Delta_{\bar{x}x} \in \mathbb{R}^{h \times \bar{n}n}$ ,  $\Lambda_{\bar{x}u} \in \mathbb{R}^{h \times \bar{n}m}$ ,  $\Lambda_{\bar{x}\bar{x}} \in \mathbb{R}^{h \times \bar{n}^2}$ . With  $\bar{\Theta} \in \mathbb{R}^{2h \times 2n\bar{n}}$  given in (29) on top of this page, we obtain

$$\bar{\Theta} \begin{bmatrix} \text{vec}(\mathbf{\Pi}_x^*) \\ \text{vec}(\mathbf{\Pi}_\phi^*) \end{bmatrix} = \begin{bmatrix} \Lambda_{\bar{x}x} \text{vec}(\mathbf{P}_k \mathbf{E}_d) \\ \Lambda_{\bar{x}x} \text{vec}(\mathbf{S}) \end{bmatrix}. \quad (30)$$

With respect to solving (30), it is formulated:

*Lemma 5:* For the given  $[0, t_E] = \cup_{i=1}^h [t_i, t_{i+1}]$ , suppose

$$\text{rank}(\Lambda_{\bar{x}x}) = n\bar{n} \quad (31)$$

and condition (14) hold, then  $\bar{\Theta}$  in (29) has full column rank.

*Proof:* We prove this by showing that

$$\bar{\Theta} \begin{bmatrix} \text{vec}(\mathbf{N}) \\ \text{vec}(\mathbf{M}) \end{bmatrix} = \mathbf{0} \quad (32)$$

necessarily implies  $\mathbf{N} = \mathbf{M} = \mathbf{0}$ . To do so, we provide an analogous analysis of  $\mathbf{x}(t)^\top \mathbf{P}_k \mathbf{N} \bar{\mathbf{x}}(t)$  and  $\mathbf{x}(t)^\top \mathbf{M} \bar{\mathbf{x}}(t)$  as in (26) and (27), respectively. According to (25), integrating the derivative over  $[t_i, t_{i+1}]$  and rearranging gives

$$\begin{aligned} & \mathbf{x}^\top(t_{i+1}) \mathbf{P}_k \mathbf{N} \bar{\mathbf{x}}(t_{i+1}) - \mathbf{x}^\top(t_i) \mathbf{P}_k \mathbf{N} \bar{\mathbf{x}}(t_i) + \int_{t_i}^{t_{i+1}} -\mathbf{x}^\top \Omega_k \mathbf{N} \bar{\mathbf{x}} \\ & - (\mathbf{B} \mathbf{u} + \mathbf{E}_d \bar{\mathbf{x}})^\top \mathbf{P}_k \mathbf{N} \bar{\mathbf{x}} + \mathbf{x}^\top \mathbf{P}_k \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{M} \bar{\mathbf{x}} \, d\tau \\ & = \int_{t_i}^{t_{i+1}} \mathbf{x}^\top \mathbf{P}_k \left( -\mathbf{A} \mathbf{N} + \mathbf{N} \bar{\mathbf{A}} + \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{M} \right) \bar{\mathbf{x}} \, d\tau. \end{aligned}$$

Corresponding to the first equation in (27), rearranging yields

$$\begin{aligned} & \mathbf{x}^\top(t_{i+1}) \mathbf{M} \bar{\mathbf{x}}(t_{i+1}) - \mathbf{x}^\top(t_i) \mathbf{M} \bar{\mathbf{x}}(t_i) \\ & + \int_{t_i}^{t_{i+1}} \mathbf{x}^\top \mathbf{C}^\top \mathbf{Q} \mathbf{C} \mathbf{N} \bar{\mathbf{x}} - (\mathbf{B} \mathbf{u} + \mathbf{E}_d \bar{\mathbf{x}})^\top \mathbf{M} \bar{\mathbf{x}} \, d\tau \\ & = \int_{t_i}^{t_{i+1}} \mathbf{x}^\top \left( \mathbf{M} \bar{\mathbf{A}} + \mathbf{A}^\top \mathbf{M} + \mathbf{C}^\top \mathbf{Q} \mathbf{C} \mathbf{N} \right) \bar{\mathbf{x}} \, d\tau. \end{aligned}$$

Choosing  $i = 1, \dots, h$ , one can observe that the  $(i)$ -th row and  $(h+i)$ -th row on the left of (32) are precisely given by the left hand side of the just stated first and second equation, respectively. Consequently, (32) is equivalent to:

$$\begin{bmatrix} \Lambda_{\bar{x}x}(\mathbf{I}_{\bar{n}} \otimes \mathbf{P}_k) \text{vec} \left( -\mathbf{A} \mathbf{N} + \mathbf{N} \bar{\mathbf{A}} + \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{M} \right) \\ \Lambda_{\bar{x}x} \text{vec} \left( \mathbf{M} \bar{\mathbf{A}} + \mathbf{A}^\top \mathbf{M} + \mathbf{C}^\top \mathbf{Q} \mathbf{C} \mathbf{N} \right) \end{bmatrix} = \mathbf{0}.$$

Due to (31) and  $\mathbf{P}_k$  invertible, it results

$$\begin{bmatrix} \mathbf{N} \\ \mathbf{M} \end{bmatrix} \bar{\mathbf{A}} = \begin{bmatrix} \mathbf{A} & -\mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \\ -\mathbf{C}^\top \mathbf{Q} \mathbf{C} & -\mathbf{A}^\top \end{bmatrix} \begin{bmatrix} \mathbf{N} \\ \mathbf{M} \end{bmatrix} \quad (33)$$

which is the homogeneous part of (13). But, condition (14) is sufficient for  $\sigma(\bar{\mathbf{A}}) \cap \sigma(\mathbf{H}) = \emptyset$  since  $\mathbf{H}$  is *Hamiltonian*. Thus, cf. [19],  $\mathbf{N} = \mathbf{M} = \mathbf{0}$  uniquely solve (33) and (32). ■

Taking the structure of  $\Lambda_{xx}$ ,  $\Lambda_{xu}$  and orthogonality of commutation matrix  $\mathbf{L}^{(n, \bar{n})}$  in (28b) into account, we conclude

*Corollary 1:* Condition (24) readily implies (31). ■

Under these conditions, the least square solution of (30) reads

$$\begin{bmatrix} \text{vec}(\mathbf{\Pi}_x^*) \\ \text{vec}(\mathbf{\Pi}_\phi^*) \end{bmatrix} = \left( \bar{\Theta}^\top \bar{\Theta} \right)^{-1} \bar{\Theta}^\top \begin{bmatrix} \Lambda_{\bar{x}x} \text{vec}(\mathbf{P}_k \mathbf{E}_d) \\ \Lambda_{\bar{x}x} \text{vec}(\mathbf{S}) \end{bmatrix}. \quad (34)$$

Following the derivations of (26) and (27), an exact solution of (30) is certainly given by the unique solution of (13). When (31) is true, this exact solution is unique and coincides with the least square solution (34); hence, (30) and (13) are equivalent sets of equations. Based on this discussion, the next theorem is self-evident again.

*Theorem 2:* For an arbitrary  $k \in \mathbb{N}_0$ , let  $\mathbf{P}_k$  satisfy (17) where  $\mathbf{K}_0$  such that  $\mathbf{A} + \alpha \mathbf{I} - \mathbf{B} \mathbf{K}_0$  is *Hurwitz* for a given  $\alpha \geq 0$ . If conditions (14) and (31) are satisfied, then  $\mathbf{\Pi}_x^*$  and  $\mathbf{\Pi}_\phi^*$  obtained from (34) uniquely solve (13). ■

*Remark 4:* Asmp. 4 ensures that  $\text{rank}(\Lambda_{\bar{x}\bar{x}}) = n\bar{n}$  can be satisfied. Hence, it must also hold in [8], [9], [17] and [18]. It is dispensable if exosystem (2a) can be excited by exploration noise:  $\bar{\mathbf{B}} \bar{\mathbf{e}}$ . This is usually possible for the part of (2a) generating  $\bar{\mathbf{y}}$  and requires minor modifications of (29).

## V. ONLINE ALGORITHM AND SIMULATION

We briefly summarize our new results by providing an online algorithm for finding  $\mathbf{u}^*(\cdot)$ . We assume that Asmp. 1-4 and condition (14) hold.

*Algorithm 1:*

- 1) *Initialization:* Give weights  $\mathbf{Q} \succeq \mathbf{0}$ ,  $\mathbf{R} \succ \mathbf{0}$  for LQTP 1 and  $\tilde{\mathbf{Q}} \succeq \mathbf{0}$ ,  $\tilde{\mathbf{R}} \succ \mathbf{0}$  for transition, i.e. LQRP 2. Find a stabilizing  $\mathbf{K}_0$  which determines a minimal degree of stability  $\bar{\alpha}_0$ , e.g. by following Remark 2. Choose a desired guaranteed degree of stability  $0 \leq \alpha < \bar{\alpha}_0$ .
- 2) *Exploration:* Run systems (2a) and (1a) excited by control input  $\mathbf{u} = -\mathbf{K}_0 \mathbf{x} + \mathbf{e}$  with exploration noise  $\mathbf{e}$  selected as, e.g., sum of sinusoidal signals (cf. Remark 3). Simultaneously, collect the measurement data (21) and (28) based on measuring  $\mathbf{x}$ ,  $\mathbf{u}$ ,  $\bar{\mathbf{x}}$  and integration over each interval  $[t_i, t_{i+1}]$ , see Section IV-A. Increment  $i \in \mathbb{N}$  until rank condition (24) holds.
- 3) *Policy-Iteration:* Iteratively calculate  $\mathbf{P}_k$  and  $\mathbf{K}_{k+1}$ ,  $k \in \mathbb{N}_0$ , from (23) unless a typical convergence criterion, e.g.  $\|\mathbf{P}_{\bar{k}} - \mathbf{P}_{\bar{k}-1}\|_F \leq \epsilon_P$ ,  $\epsilon_P \in \mathbb{R}^{>0}$ , is satisfied at  $k = \bar{k}$ . Under premise that  $\bar{k}$  is sufficiently large, we have  $\mathbf{K}_{\bar{k}+1} \approx \mathbf{K}^*$  based on Theorem 1.

For an arbitrary  $0 \leq k \leq \bar{k}$ , we exactly obtain  $\mathbf{B} = \mathbf{P}_k^{-1} \mathbf{K}_{k+1}^\top \mathbf{R}$  and  $\mathbf{E}_d$  from 3<sup>rd</sup> argument of (23).

4) *One-Step Calculation*: For an arbitrary  $0 \leq k \leq \bar{k}$ , calculate  $\mathbf{\Pi}_x^*$  and  $\mathbf{\Pi}_\phi^*$  from (34), which *exactly* solve (13) based on Theorem 2, and  $\mathbf{F}^* = -\mathbf{R}^{-1} \mathbf{B}^\top \mathbf{\Pi}_\phi^*$ .

5) *Implementation*: To achieve C2-4) on  $[t_E, \infty)$ , cf. Section I, apply  $\mathbf{u}^* \approx -\mathbf{K}_{\bar{k}+1}^-(\mathbf{x} - \mathbf{\Pi}_x^* \bar{\mathbf{x}}) + \mathbf{F}^* \bar{\mathbf{x}}$ .

*Remark 5*: As in [8], [9], [12], [17], [18], [20] and [21], we have to assume that all measurements are *not* corrupted by *unknown* sensor noise. In general, it would be of interest to make these methods more robust in this aspect.

In the succeeding simulation example of an over-actuated system, i.e.  $\text{rank}(\mathbf{B}) > \text{rank}(\mathbf{C})$ , we implement Algorithm 1 and emphasize our contribution C4) by two observations. First, if output regulation is feasible, our approach is suited for *almost exact tracking*, cf. Section I. Second, the additional actuators are used more efficiently to save input-energy than in [8], [9] whereas our calculation procedure is less complicated at the same time.

*Example*: We consider a non-minimum phase and over-actuated LTI-system with  $(\mathbf{C}^\top, \mathbf{A}, \mathbf{B}, \mathbf{E}_d, \mathbf{D}_d) :=$

$$\left( \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 & 0 \\ 2 & 1 & 1 \\ 0 & 0 & 2 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 0 & -2 & 0 \\ 0 & 0 & 0 & 3 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{0} \right).$$

The exosystem  $(\bar{\mathbf{C}}, \bar{\mathbf{A}})$  and initial value  $\bar{\mathbf{x}}_0$  are given such that  $\bar{\mathbf{x}}(t)^\top = [t \ 1 \ \cos(t)/4 \ \sin(t)/4]$  and scalar reference  $\bar{y}(t) = (t-5)/4$  (---) result, i.e. we regard a polynomial reference and sinusoidal disturbances.

In the sequel, the data (4) of the dynamics is assumed to be **unknown**. Following Remark 2, we may obtain a stabilizing

$$\mathbf{K}_0 = \begin{bmatrix} 9.56 & 7.1 & 0.96 \\ 0.99 & 0.96 & 6.7 \end{bmatrix}.$$

Without knowing  $\mathbf{A}$  and  $\mathbf{B}$ , diagonal weights such as  $\tilde{\mathbf{Q}} = \text{diag}(10, 10, 5)$  and  $\tilde{\mathbf{R}} = \mathbf{I}_2$  seem appropriate with respect to LQRP 2. Here, the first two entries of  $\tilde{\mathbf{Q}}$  are equally chosen since both, 1<sup>st</sup> state  $x_1$  and 2<sup>nd</sup> state  $x_2$ , contribute similarly to  $y$  based on the known  $\mathbf{C}$ . We cannot anticipate, however, that we penalized  $x_1$  and its derivative  $\dot{x}_1 = x_2$ . This causes a counterbalance for the chosen  $\tilde{\mathbf{Q}}$  such that the real part of a closed-loop eigenvalue is bounded from the left by  $-1$  for any  $\mathbf{K}_k$ ,  $k \geq 1$ . In fact, this is even true when we scale  $\tilde{\mathbf{Q}}$  by an arbitrarily large positive scalar. Clearly, this indicates how useful our extension of the method by [12] can be as it guarantees in general that all closed-loop eigenvalues lie to the left of a defined  $-\alpha$  independent of  $\tilde{\mathbf{Q}}$ ,  $\mathbf{R}$ . In particular, we might know that  $\alpha < 1.5$  satisfies the requirements in Lemma 3, e.g. by a study as in Remark 2. Indeed, we have  $\max_{\lambda \in \sigma(\mathbf{A} - \mathbf{B}\mathbf{K}_0)} \text{Re}(\lambda) < -1.7$ . By choosing  $\alpha = 1$ , all closed-loop eigenvalues lie certainly to the left of  $-1$ .

In view of tracking  $\bar{y}(t)$  (---), we consider  $i = 1, 2$  different sets of weights:  $\mathbf{Q}_1 = 200$ ,  $\mathbf{R}_1 = 1/10 \cdot \mathbf{I}_2$  (—) and  $\mathbf{Q}_2 = 0.5$ ,  $\mathbf{R}_2 = \mathbf{R}_1$  (—). In the first case, we desire to track  $\bar{y}(t)$  closely while in the second tracking errors are tolerated to save input energy. Since the data (4) is not available, for neither of the cases  $i = 1, 2$  we are able to

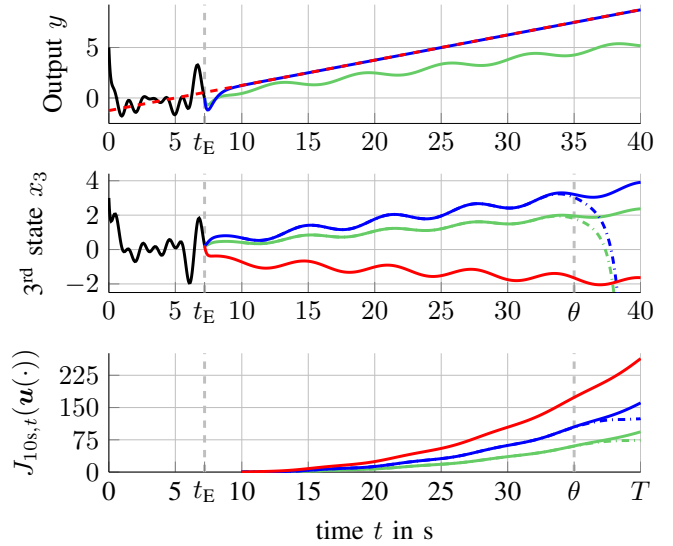


Fig. 1. Simulation results with exploration  $[0, t_E]$  (—) and application  $[t_E, T]$  for two different weights  $\mathbf{Q}_1 = 200$  (—) and  $\mathbf{Q}_2 = 0.5$  (—) for tracking of (---). Comparison with each finite horizon  $T = 40$  s optimal solution (---), (---) and “optimal” output regulation (—) of [8], [9].

exactly calculate  $\mathbf{u}_i^*(\cdot)$  as in Lemma 1. Clearly, this is in general severe considering 3) of Lemma 1. However, we may apply our developed Algorithm 1. The information required for *initialization* 1) has already been provided above. For 2), the *exploration* noise is defined as  $e = \sum_{\omega=2}^6 \mathbf{r}_\omega \sin(\frac{\omega}{s} t)$  with randomly chosen  $\mathbf{r}_\omega \in \mathbb{R}^2$ ,  $\mathbf{r}_\omega \neq \mathbf{0}$ . This excitation by 5 distinct frequencies is *sufficiently rich*, cf. Remark 3. Hence, condition (24) with  $\frac{1}{2}n(n+1) + nm + n\bar{n} = 24$  holds after  $t_E = 24 \cdot 0.3 \text{ s} = 7.2 \text{ s}$  of exploration for  $t_{i+1} - t_i = 0.3 \text{ s}$ . The exploration phase (—) is illustrated in Fig. 1 for output  $y$  and the 3<sup>rd</sup> state  $x_3$ , where  $\mathbf{x}_0^\top = [0 \ 5 \ 3]$  and  $\mathbf{x}(t_E)^\top = [0.089 \ 0.092 \ 0.232]$ .

Subsequently, the *policy iteration* 3) is performed. Reminding Section I, [17] and [18] carry out this step for the augmented  $n + \bar{n}$ -dimensional system in order to obtain  $\mathbf{u}_{\text{aug}} = -\mathbf{R}^{-1} [\mathbf{B}^\top \ 0] \mathbf{P}_{\text{aug},k} [\mathbf{x}^\top \ \bar{\mathbf{x}}^\top]^\top$ . This also determines the tracking performance which obviously depends on the accuracy of  $\mathbf{P}_{\text{aug},k}$ . In contrast, our tracking performance does not depend on the accuracy of  $\mathbf{P}_k$ . Hence,  $\epsilon_P = 0.05$  serves our needs for which we obtain  $\|\mathbf{P}_{\bar{k}} - \mathbf{P}\|_F = 0.003$  and  $\|\mathbf{K}_{\bar{k}+1} - \mathbf{K}^*\|_F = 0.0006$  after only  $\bar{k} = 3$  steps. By 3) and our *one-step calculation* 4) for  $i = 1, 2$ , we calculate  $\mathbf{B}$ ,  $\mathbf{E}_d$  and  $\mathbf{\Pi}_{x,i}^*$ ,  $\mathbf{\Pi}_{\phi,i}^*$ . The Frobenius norm of each deviation from its true value is smaller than  $3.6 \cdot 10^{-12}$  and, thus, lies in range of the numerical precision of a calculation by (13).

In summary, though not having knowledge of (4) we have found  $\mathbf{u}^*(\cdot)$  such that the specifications of an optimal transition by LQRP 2 are approximatively met with a guaranteed degree of stability and these of stationarily optimal tracking by LQTP 1 with unbounded cost are exactly met. That is, the desired properties in Lemma 1 and 2 are achieved.

The phase of application  $[t_E, T]$  with  $T = 40$  s is displayed in Fig. 1. Obviously,  $\mathbf{u}_2^*(\cdot)$  (—) with  $\mathbf{Q}_2 = 0.5$  tracks the reference  $\bar{y}$  (---) less closely and attenuates disturbances

less than  $\mathbf{u}_1^*(\cdot)$  (—) with  $Q_1 = 200$ . Thus,  $\mathbf{u}_2^*(\cdot)$  saves 58.6% input-energy compared with  $\mathbf{u}_1^*(\cdot)$ . But on the other hand,  $\mathbf{u}_1^*(\cdot)$  leads to almost exact tracking with stationary tracking error only about 1‰. When the transition is over:  $\tilde{\mathbf{x}}(t) \approx \mathbf{0}$  for  $t \geq 10$  s, we are interested in the stationary cost  $J_{10s,t}(\mathbf{u}(\cdot))$ . For comparison, we consider exact tracking by “optimal” output regulation (OR) given in [8], [9]. There, the stationary state  $\mathbf{\Pi}_{x,OR}\bar{\mathbf{x}}$  and control  $\mathbf{u}_{OR}(\cdot) = \mathbf{\Gamma}_{OR}\bar{\mathbf{x}}$  (—) are the solution of

$$\begin{aligned} & \min_{\mathbf{\Pi}_x, \mathbf{\Gamma}} \text{trace} \left( \mathbf{\Gamma}^T \mathbf{R}_1 \mathbf{\Gamma} \right) \\ \text{s.t.: } & \mathbf{\Pi}_x \bar{\mathbf{A}} = \mathbf{A} \mathbf{\Pi}_x + \mathbf{B} \mathbf{\Gamma} + \mathbf{E}_d \text{ and } \mathbf{C} \mathbf{\Pi}_x + \mathbf{D}_d = \bar{\mathbf{C}}. \end{aligned}$$

This optimization problem (OP) was introduced by [14] in case that the solution of the regulator equations in [19], i.e. the constraints, is not unique due to over-actuation. It intends to provide an optimized solution in view of a cost functional such as (3). But since the OP is generally not equivalent it results in sub-optimal solutions as already indicated in [4].

As a consequence, while  $\mathbf{u}_{OR}(\cdot)$  (—) and  $\mathbf{u}_1^*(\cdot)$  (—) lead to an equivalent tracking performance  $\bar{y}^*(t) \approx \bar{y}_{OR}(t)$ ,  $\mathbf{u}_{OR}(\cdot)$  requires stationarily 64% more input energy than  $\mathbf{u}_1^*(\cdot)$ ; note that cost  $J_{10s,t}(\mathbf{u}(\cdot))$  equals the input-energy if  $y(t) - \bar{y}(t) \approx 0$ . Indeed,  $\mathbf{u}_1^*(\cdot)$  outperforms  $\mathbf{u}_{OR}(\cdot)$  which is approved by the observations in [4]. This is reasoned as follows. Most of the entries in  $\mathbf{\Pi}_x^* = [\pi_{x,i,j}^*]$  and  $\mathbf{\Pi}_{x,OR} = [\pi_{x,OR,i,j}]$  approximatively coincide. More precisely, we have indeed  $\lim_{Q_1 \rightarrow \infty} \pi_{x,i,j}^* = \pi_{x,OR,i,j}$ , besides  $\pi_{x,3,1}^* \not\approx \pi_{x,OR,3,1}$  and  $\pi_{x,3,2}^* \not\approx \pi_{x,OR,3,2}$ . Hence, it results an essentially different stationary behavior of 3<sup>rd</sup> state  $x_3$  for  $\mathbf{u}_1^*(\cdot)$  and  $\mathbf{u}_{OR}(\cdot)$ , see Fig. 1. Under the exact tracking constraint, this reflects the degrees of freedom available due to the over-actuation. These are clearly more efficiently exploited by  $\mathbf{u}_1^*(\cdot)$ . In addition, to obtain  $\mathbf{u}_{OR}(\cdot)$  in [8], [9], the optimization problem above has to be solved online whereas we only need to solve a single set of equations (30) in step 4) to derive  $\mathbf{F}^*$ .

In view of Lemma 2 with  $T_1 = 10$  s, we compare the stationarily agreeable plans  $\mathbf{u}_1^*(\cdot)$  (—) and  $\mathbf{u}_2^*(\cdot)$  (—) with their corresponding finite horizon  $[T_1, 40$  s] optimal solutions  $\mathbf{u}_{T,1}^*(\cdot)$  (---) and  $\mathbf{u}_{T,2}^*(\cdot)$  (---), respectively. As we see from Fig. 1, we can choose  $\theta = 35$  s to guarantee a relative cost increase of less than 1% for  $i = 1, 2$  on  $[10$  s,  $35$  s]. Switching from  $\mathbf{u}_i^*(\cdot)$  to  $\mathbf{u}_{T,i}^*(\cdot)$  on  $[35$  s,  $40$  s] would allow us to save input-energy. However, the unstable characteristics of  $x_3$  for  $\mathbf{u}_{T,i}^*(\cdot)$  on  $[\theta, T]$  are unacceptable. To prevent this, one typically considers a quadratic end-cost for penalizing  $|x_3(T)|$  in (3). But then, the mentioned performance gain of  $\mathbf{u}_{T,i}^*(\cdot)$  on  $[\theta, T]$  drops significantly and we may simply stick to  $\mathbf{u}_i^*(\cdot)$  which is considerably easier to implement.

## VI. CONCLUSION

For the first time, we derived a method to determine an LQ optimal tracking control for problems with unbounded cost without the need to know the system dynamics in advance. The control consists of two independent parts. A state feedback guarantees an optimal transition with specified

degree of stability. It solves an LQR problem with individual weights and is obtained by extending the iterative ADP approach in [9], [12]. The second part is a static pre-filter resulting in optimal stationary tracking and is provided by a new one-step calculation. It adequately solves an infinite horizon LQT problem with generally unbounded cost and approximates finite horizon optimal solutions. For over-actuated systems, we showed that the new approach leads to a less complicated calculation and a more efficient operation in comparison with [8], [9].

## REFERENCES

- [1] B. D. O. Anderson and J. B. Moore, *Optimal Control: Linear Quadratic Methods*. Dover Publications, Inc., 2007.
- [2] Z. Artstein and A. Leizarowitz, “Tracking periodic signals with the overtaking criterion,” *IEEE Transactions on Automatic Control*, vol. 30, no. 11, pp. 1123–1126, 1985.
- [3] S. Bernhard, “Time-invariant control in LQ optimal tracking: An alternative to output regulation,” *IFAC-PapersOnLine*, vol. 50, no. 1, pp. 4912 – 4919, 2017, 20th IFAC World Congress.
- [4] S. Bernhard and J. Adamy, “Static optimal decoupling control for linear over-actuated systems regarding time-varying references,” in *2017 American Control Conf. (ACC)*, 2017, pp. 1049–1055.
- [5] S. Boyd, L. E. Ghaoui, E. Feron, and V. Balakrishnan, *Linear Matrix Inequalities in System and Control Theory*. SIAM, 1994.
- [6] D. A. Carlson, “Uniformly overtaking and weakly overtaking optimal solutions in infinite-horizon optimal control: When optimal solutions are agreeable,” *Journal of Optimization Theory and Applications*, vol. 64, no. 1, pp. 55–69, 1990.
- [7] D. A. Carlson and A. Haurie, *Infinite Horizon Optimal Control*. Springer-Verlag Berlin Heidelberg, 1987.
- [8] W. Gao and Z.-P. Jiang, “Linear optimal tracking control: An adaptive dynamic programming approach,” in *2015 American Control Conference (ACC)*, 2015, pp. 4929–4934.
- [9] —, “Adaptive dynamic programming approach and adaptive optimal output regulation of linear systems,” *IEEE Transactions on Automatic Control*, vol. 61, no. 12, pp. 4164–4169, 2016.
- [10] H. Halkin, “Necessary conditions for optimal control problems with infinite horizons,” *Econometrica*, vol. 42, no. 2, pp. 267–272, 1974.
- [11] P. Ioannou and J. Sun, *Robust Adaptive Control*. Prentice Hall, 1996.
- [12] Y. Jiang and Z.-P. Jiang, “Computational adaptive optimal control for continuous-time linear systems with completely unknown dynamics,” *Automatica*, vol. 48, no. 10, pp. 2699–2704, 2012.
- [13] D. L. Kleinman, “On an iterative technique for Riccati equation computations,” *IEEE Transactions on Automatic Control*, vol. 13, no. 1, pp. 114–115, 1968.
- [14] A. J. Krener, “The construction of optimal linear and nonlinear regulators,” *Systems, Models and Feedback: Theory and Applications*, vol. 12, pp. 301–322, 1992.
- [15] F. L. Lewis and D. Vrabie, “Reinforcement learning and adaptive dynamic programming for feedback control,” *IEEE Control Systems Magazine*, vol. 9, no. 3, pp. 32–50, 2009.
- [16] F. L. Lewis, D. Vrabie, and K. G. Vamvoudakis, “Reinforcement learning and feedback control,” *IEEE Control Systems Magazine*, vol. 32, no. 6, pp. 76–105, 2012.
- [17] H. Modares and F. L. Lewis, “Linear quadratic tracking control of partially-unknown continuous-time systems using reinforcement learning,” *IEEE Transactions on Automatic Control*, vol. 59, no. 11, pp. 3051–3056, 2014.
- [18] C. Qin, H. Zhang, and Y. Luo, “Online optimal tracking control of continuous-time linear systems with unknown dynamics by using adaptive dynamic programming,” *International Journal of Control*, vol. 87, no. 5, pp. 1000–1009, 2014.
- [19] H. Trentelman, A. A. Stoorvogel, and M. Hautus, *Control Theory for Linear Systems*. Springer-Verlag London, 2001.
- [20] K. G. Vamvoudakis, “Optimal trajectory output tracking control with a Q-learning algorithm,” in *2016 American Control Conference (ACC)*, 2016, pp. 5752–5757.
- [21] D. Vrabie, O. Pastravanu, M. Abu-Khalaf, and F. L. Lewis, “Adaptive optimal control for continuous-time linear-systems based on policy iteration,” *Automatica*, vol. 45, no. 2, pp. 477–484, 2009.