# KULLBACK-LEIBLER-QUADRATIC OPTIMAL CONTROL\*

NEIL CAMMARDELLA<sup>†</sup>, ANA BUŠIĆ<sup>‡</sup>, AND SEAN P. MEYN<sup>§</sup>

Abstract. This paper presents approaches to mean-field control, motivated by distributed control of multiagent systems. Control solutions are based on a convex optimization problem, whose domain is a convex set of probability mass functions (pmfs). The main contributions follow: (1) Kullback-Leibler-quadratic (KLQ) optimal control is a special case in which the objective function is composed of a control cost in the form of Kullback-Leibler divergence between a candidate pmf and the nominal, plus a quadratic cost on the sequence of marginals. Theory in this paper extends prior work on deterministic control systems, establishing that the optimal solution is an exponential tilting of the nominal pmf. Transform techniques are introduced to reduce complexity of the KLQ solution, motivated by the need to consider time horizons that are much longer than the intersampling times required for reliable control. (2) Infinite-horizon KLQ leads to a state feedback control solution with attractive properties. It can be expressed as state feedback, in which the state is the sequence of marginal pmfs, or an open loop solution is obtained that is more easily computed. (3) Numerical experiments are surveyed in an application of distributed control of residential loads to provide grid services, similar to utility-scale battery storage. The results show that KLQ optimal control enables the aggregate power consumption of a collection of flexible loads to track a time-varying reference signal, while simultaneously ensuring each individual load satisfies its own quality of service

Key words. mean-field games, distributed control, demand dispatch

MSC codes. 90C40, 93E20, 90C46, 93E35, 60J20

**DOI.** 10.1137/21M1433654

- 1. Introduction. The goal of this paper is to obtain control solutions for mean-field models. The optimization problems considered are generalizations of standard Markov decision process (MDP) objectives, in both finite-horizon and average-cost settings.
- **1.1. Mean-field control.** The mean-field control problem is an approach to distributed control of a collection of  $\mathcal N$  homogeneous "agents," with  $\mathcal N\gg 1$ , modeled as discrete-time stochastic systems, with state processes at time k denoted  $\{X_k^i:1\leq i\leq \mathcal N\}$ . To avoid a long detour on notation it is assumed that the common state space  $\mathsf X$  is finite. For a single value k and time horizon  $K\geq 1$ , the empirical distributions are denoted

<sup>\*</sup>Received by the editors July 14, 2021; accepted for publication (in revised form) July 11, 2023; published electronically October 25, 2023.

https://doi.org/10.1137/21M1433654

Funding: The authors received support from National Science Foundation grant EPCN 1935389 and from the France 2030 Plan, PEPR TASE, AI-NRGY project ANR-22-PETA-0044.

 $<sup>^{\</sup>dagger}$ Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL 32611 USA (ncammardella@ufl.edu).

<sup>&</sup>lt;sup>‡</sup>Inria and DI ENS, École Normale Supérieure, CNRS, PSL Research University, Paris, France (ana.busic@inria.fr, https://www.di.ens.fr/~busic/).

<sup>§</sup>Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL 32611 USA, and Inria International Chair, Paris, France (meyn@ece.ufl.edu, http://www.meyn.ece.ufl.edu/).

(1.1a) 
$$p^{\mathcal{N}}(\vec{x}) = \frac{1}{\mathcal{N}} \sum_{i=1}^{\mathcal{N}} \mathbb{I}\{(X_0^i, \dots, X_K^i) = \vec{x}\}, \qquad \vec{x} \in \mathsf{X}^{K+1},$$

$$(1.1b) \qquad \quad \nu_k^{\mathcal{N}}(x) = \frac{1}{\mathcal{N}} \sum_{i=1}^{\mathcal{N}} \mathbb{I}\{X_k^i = x\}\,, \qquad \qquad x \in \mathsf{X}\,,$$

where  $\vec{x} = (x_0, \dots, x_K)$  denotes an arbitrary element of  $X^{K+1}$ . The set of probability mass functions (pmfs) on  $X^{K+1}$  is denoted by  $S(X^{K+1})$  for  $K \ge 1$  and S(X) for K = 0.

The integer  $\mathcal{N}$  is regarded as a parameter in mean-field theory, and assumptions imply that there is convergence as  $\mathcal{N} \to \infty$ ,

$$\lim_{\mathcal{N} \to \infty} p^{\mathcal{N}}(\vec{x}) = p(\vec{x}) \,, \qquad \lim_{\mathcal{N} \to \infty} \nu_k^{\mathcal{N}}(x_k) = \nu_k(x) \,,$$

where  $\nu_k \in \mathcal{S}(X)$  is the kth marginal of  $p \in \mathcal{S}(X^{K+1})$  for  $0 \le k \le K$ .

In this paper this limit is achieved by assuming homogeneity of the statistics of each agent: for each i the state evolution is consistent with p:

(1.2) 
$$P\{X_{k+1}^i = x_{k+1} \mid (X_0^i, \dots, X_k^i) = \vec{x}_0^k\} = p(x_{k+1} \mid \vec{x}_0^k),$$

where the conditional pmfs are obtained from the Bayes rule.

The paper concerns design of p to balance two objectives, based on a reference signal  $\{r_k\}$ , and function  $\mathcal{Y} \colon \mathsf{X} \to \mathbb{R}$ :

- (i)  $\nu_k \sim \nu_k^0$ , where  $\{\nu_k^0\}$  models nominal behavior.
- (ii)  $\langle \nu_k, \mathcal{Y} \rangle := \sum_{x \in \mathsf{X}} \nu_k(x) \mathcal{Y}(x) \approx r_k$ .

The agents considered in section 4 represent a population of residential water heaters, and  $\mathcal{Y}: \mathsf{X} \to \mathbb{R}_+$  is chosen so that  $\langle \nu_k^{\mathcal{N}}, \mathcal{Y} \rangle$  is the average power consumption over the population of loads.

Two approaches to design are developed in this paper.

**Feedforward control.** A sequence  $\{C_k : 1 \le k \le K\}$  of real-valued cost functions on the marginals is specified, and  $p^*$  is obtained as the solution to

(1.3) 
$$J^{\star}(\nu_0^0) = \min_{p} \sum_{k=1}^{K} C_k(\nu_k),$$

where the minimum is over all pmfs with first marginal  $\nu_0^0$ . The two goals motivate the objective function

(1.4) 
$$C_k(\nu) = \mathcal{D}(\nu, \nu_k^0) + \frac{\kappa}{2} \left[ \langle \nu, \mathcal{Y} \rangle - r_k \right]^2, \qquad \nu \in \mathcal{S}(\mathsf{X}),$$

in which  $\kappa > 0$  is a penalty parameter, and  $\mathcal{D}$  penaltzes deviation from nominal behavior. The finite-horizon optimal control problem is thus

(1.5) 
$$J^{\star}(\nu_0^0) = \min_{p} \sum_{k=1}^{K} \left[ \mathcal{D}(\nu_k, \nu_k^0) + \frac{\kappa}{2} \left[ \langle \nu_k, \mathcal{Y} \rangle - r_k \right]^2 \right].$$

It is envisioned that this finite horizon optimal control problem will be a component of a model predictive control (MPC) strategy, with time horizons for computation updates dictated by performance requirements and model accuracy.

**Feedback control.** If the nominal model is Markovian, then the evolution of the marginals follows the dynamics of a controlled nonlinear state space model,

(1.6) 
$$\nu_{k+1} = f_k(\nu_k, \phi_k), \quad k \ge 0, \quad \nu_0^0 \text{ given},$$

where  $\{\phi_k\}$  is the input sequence, evolving on an abstract set  $\Phi$ . A feedback policy takes the form  $\phi_k = \mathcal{K}_k(\nu_k)$ .

Design choices for  $\mathcal{K}_k$  are proposed based on an infinite-horizon solution of (1.5). Justification requires further assumptions, including time-homogeneous dynamics for (1.6), which holds if the nominal model is a time-homogeneous Markov chain.

1.2. MDPs and mean-field control. The Markovian assumption for the nominal model is based on the standard controlled Markov chain model used in MDPs.

The model considered here is specified by a state space denoted S and input space U, and we denote  $X := S \times U$  (assumed finite). The joint state-input process is denoted  $X = \{X_k = (S_k, U_k) : k \ge 0\}$ . In finite-horizon optimal control the model includes a sequence of controlled transition matrices  $\{T_k : k \ge 0\}$  and cost functions  $\{c_k : k \ge 0\}$  with  $c_k : X \to \mathbb{R}$  for each k.

The dynamics of  $\mathbf{x} = (\mathbf{S}, \mathbf{U}) = \{S_k, U_k : k \geq 0\}$  are determined by the transition matrices as follows. It is assumed that  $\mathbf{X}$  is adapted to a filtration  $\{\mathcal{F}_k : k \geq 0\}$  (so that  $X_k$  is  $\mathcal{F}_k$ -measurable for each k), and

(1.7) 
$$P\{S_{k+1} = s' \mid \mathcal{F}_k; \ S_k = s, \ U_k = u\} = T_k(x, s'), \qquad x = (s, u) \in X, s' \in S.$$

The set of functions from S to the simplex  $\mathcal{S}(\mathsf{U})$  is denoted  $\Phi$ , and we let  $\phi$  denote a generic element of  $\Phi$ . The decision rule defining the input sequence is assumed to be Markovian: with  $\phi_k \in \Phi$  for each k,

(1.8) 
$$P\{U_k = u \mid \mathcal{F}_{k-1}; \ S_k = s\} = \phi_k(u \mid s), \qquad x = (s, u) \in X.$$

The finite-horizon optimal control problem of MDP theory is a special case of (1.3), in which  $C_k$  linear for each k; in this case  $C_k(\nu_k) = \langle \nu_k, c_k \rangle = \sum_{x \in \mathsf{X}} \nu_k(x) c_k(x)$  for each k, and the sum on the right-hand side of (1.3) may be expressed

$$\sum_{k=1}^K \langle \nu_k, c_k \rangle = \sum_{k=1}^K \mathsf{E}[c_k(X_k)] \,, \qquad X_k \sim \nu_k \,,$$

where X evolves according to the controlled Markovian dynamics. This interpretation is the first step in the linear programming (LP) approach to MDPs introduced by Manne [5, 34]. The second step is to recognize that the dynamics can be expressed as a sequence of linear constraints on the marginals,

$$(1.9) \qquad \sum_{u'} \nu_k(s',u') = \sum_{s,u} \nu_{k-1}(s,u) T_{k-1}(x,s') \,, \quad s' \in \mathsf{S} \,, \ 1 \leq k \leq K \,, \quad \nu_0^0 \text{ given}.$$

Another special case of (1.3) is variance-penalized optimal control, for which  $C_k(\nu_k) = \langle \nu_k, c \rangle + \kappa [\langle \nu_k, c^2 \rangle - \langle \nu_k, c \rangle^2]$ , with  $\kappa > 0$  a penalty parameter. The solution to the optimization problem (1.3) can be expressed using a randomized state feedback policy of the form (1.8) [2, 41, 36].

1.3. Kullback-Leibler-quadratic control. In this approach to feedforward control we choose a Markovian model of the form (1.7), (1.8) to define nominal behavior: for a collection  $\{\phi_k^0\} \subset \Phi$ ,

(1.10a) 
$$p^{0}(\vec{x}) = \nu_{0}^{0}(x_{0})P_{0}^{0}(x_{0}, x_{1})P_{1}^{0}(x_{1}, x_{2})\cdots P_{K-1}^{0}(x_{K-1}, x_{K}),$$

(1.10b) 
$$P_k^0(x, x') = T_k(x, s') \phi_{k+1}^0(u' \mid s'), \qquad x, x' \in X.$$

Any other  $\{\phi_k\} \subset \Phi$  defines a Markov chain X with transition matrices,

$$(1.11) P_k(x,x') := P\{X_{k+1} = x' \mid X_k = x\} = T_k(x,s') \phi_{k+1}(u' \mid s').$$

The marginals evolve according to linear dynamics, similar to (1.9),

$$(1.12) \nu_k = \nu_{k-1} P_{k-1} \,, 1 \le k \le K,$$

in which  $\nu_k$  is interpreted as an *n*-dimensional row vector with  $n = |\mathsf{X}|$ .

We obtain a convex program by optimizing over  $\{\nu_k\}$ , similar to the LP approach of [34]. Scalar variables  $\{\gamma_k\}$  are introduced to simplify the objective, in anticipation of a Lagrangian decomposition:

(1.13a) 
$$J^{\star}(\nu_0^0) := \min_{\nu, \gamma} \left[ \sum_{k=1}^K \mathcal{D}(\nu_k, \nu_k^0) + \frac{\kappa}{2} \sum_{k=1}^K \gamma_k^2 \right]$$

(1.13b) s.t. 
$$\gamma_k = \langle \nu_k, \mathcal{Y} \rangle - r_k$$
,

(1.13c) 
$$\sum_{u'} \nu_k(s', u') = \sum_{s, u} \nu_{k-1}(s, u) T_{k-1}(x, s'), \qquad 1 \le k \le K.$$

The relative entropy rate is adopted as the cost of deviation:

(1.14) 
$$\mathcal{D}(\nu_k, \nu_k^0) := \sum_{s,u} \nu_k(s, u) \log \left( \frac{\phi_k(u \mid s)}{\phi_k^0(u \mid s)} \right).$$

The terminology is justified through the following steps. First, we have seen that any randomized policy gives rise to a pmf  $p \in \mathcal{S}(X^{K+1})$  that is Markovian:  $p(\vec{x}) = \nu_0^0(x_0)P_0(x_0,x_1)P_1(x_1,x_2)\cdots P_{K-1}(x_{K-1},x_K)$ . The relative entropy (also known as Kullback–Leibler divergence) is the mean log-likelihood:

(1.15) 
$$D(p||p^0) = \sum L(\vec{x}) p(\vec{x}),$$

where  $L = \log(p/p^0)$  is an extended-real-valued function on  $\mathsf{X}^{K+1}$ . The expression for  $P_k$  in (1.11) and the analogous formula for  $P_k^0$  using  $\varphi_{k+1}^0$  gives

$$(1.16) \quad L(\vec{x}) = \log\left(\frac{p(\vec{x})}{p^0(\vec{x})}\right) = \sum_{k=0}^{K-1} \log\left(\frac{P_k(x_k, x_{k+1})}{P_k^0(x_k, x_{k+1})}\right) = \sum_{k=1}^K \log\left(\frac{\varphi_k(u_k \mid s_k)}{\varphi_k^0(u_k \mid s_k)}\right).$$

Consequently,  $D(p||p^0) = \sum_{k=1}^K \mathcal{D}(\nu_k, \nu_k^0)$ .

The simple proof of Proposition 1.1 may be found in [13].

PROPOSITION 1.1. With  $\mathcal{D}$  chosen as the relative entropy rate (1.14), the optimization problem (1.13) is convex in  $\{\nu_k, \gamma_k : 1 \leq k \leq K\}$ . Furthermore, the linear constraints in (1.13c) are equivalent to (1.12).

1.4. Motivation from linear systems theory. The approach to feedback control proposed in section 3 begins with consideration of the infinite-horizon KLQ problem. This is tractable only subject to additional assumptions.

It is assumed that the nominal model is a time-homogeneous Markov chain and that the reference signal is *constant*,  $r_k \equiv r$ ,  $k \geq 0$ . On optimizing for each  $r \in \mathbb{R}$  we obtain a continuous family of optimizers,  $\{ \varphi_k^{\star}(u \mid s; r) : (s, u) \in \mathsf{X}, k \geq 0, \ r \in \mathbb{R} \}$ . A potentially useful policy for tracking is then

$$\phi_k(u \mid s) = \phi_k^{\star}(u \mid s; r_k), \qquad (s, u) \in X, \ k \ge 0.$$

Motivation for this approach may be found in the theory of optimal control for linear systems.

Consider the linear system with n-dimensional state X, m-dimensional input U, and scalar output Y, evolving as

$$(1.18) X_{k+1} = AX_k + BU_k + N_{k+1}, Y_k = C^T X_k + W_{k+1},$$

in which  $\{N_{k+1}, W_{k+1}\}$  are independent and identically distributed (i.i.d.), mutually independent, with zero mean and finite covariances. The cost is quadratic,  $c(x, u; r) = (y - r)^2 + u^T R u$  with R > 0.

The goal is to solve the average-cost optimal control problem. The solution is obtained via state-augmentation: define  $X_k^r = [X_k; r_k]$ , where  $r_{k+1} = r_k = r$  defines the dynamics. The solution is linear state feedback,

$$(1.19) U_k = -K^* X_k + G^* r, k \ge 0,$$

where  $[K^*; G^*]$  is the optimal gain. The optimal gain does *not* depend on r or the distribution of  $N_k$  or  $W_k$ .

The special case in which the disturbances are *zero* is most closely related to the nonlinear control problem considered in section 3. Consider the objective

$$J_K^{\star}(x) = \min_{U} \sum_{n=0}^{K} c(X_k, U_k) = \min_{u} \left\{ c(x, u; r) + J_{K-1}^{\star}(Ax + Bu) \right\}, \qquad X_0 = x \in \mathbb{R}^n.$$

It is not useful to let  $K \to \infty$  without modification, since the cost c(x, u; r) is never zero. This is why the relative value functions  $h_K(x) = J_K^{\star}(x) - J_K^{\star}(0)$  are introduced, which solve the Bellman equation in modified form,

$$\eta_K + h_K(x) = \min_u \{ c(x, u; r) + h_{K-1}(Ax + Bu) \}, \qquad X_0 = x \in \mathbb{R}^n,$$

with  $\eta_K = J_K^{\star}(0) - J_{K-1}^{\star}(0)$ . As  $K \to \infty$ , the pair  $(\eta_K, h_K)$  converge to a solution to the average-cost optimality equation (ACOE),

$$\eta^* + h^*(x) = \min_{u} \left\{ c(x, u; r) + h^*(Ax + Bu) \right\}, \qquad x \in \mathbb{R}^n,$$

whose minimizer is precisely (1.19). The proof is standard, though usually presented in the purely stochastic setting. It is especially simple in this LQR setting since each of the functions  $\{h_N\}$  is quadratic [41, 36].

When r is time varying, it is standard practice to apply the "hack"

$$(1.20) U_k = -K^* X_k + G^* r_k, k \ge 0.$$

The most compelling motivation is found in the deterministic, continuous time setting: under mild conditions, the return difference equation tells us that the closed loop dynamics from reference input to output are *passive* [3]. Passivity is lost for discrete time models but can be expected to hold approximately when the discrete time model is obtained from sampling a continuous time system.

- **1.5.** Main results. The contributions of this paper fall into three categories:
- (1) Feedforward control. Consideration of the dual of the convex optimization problem (1.13) leads to many insights. The main conclusions summarized here are a special case of Theorem 2.1.

THEOREM 1.2. [KLQ solution]. Consider the convex program (1.13). An optimizer  $\{\phi_k^*: 1 \le k \le K\}$  exists, is unique, and is of the form

(1.21a) 
$$\phi_k^{\star}(u \mid s) = \phi_k^{0}(u \mid s) \exp\left(\sum_{s'} T_k(x, s') g_{k+1}^{\star}(s') + \lambda_k^{\star} \mathcal{Y}(s, u) - g_k^{\star}(s)\right),$$

$$(1.21b) where g_k^{\star}(s) = \log \left( \sum_{u} \phi_k^0(u \mid s) \exp \left( \sum_{s'} T_k(x, s') g_{k+1}^{\star}(s') + \lambda_k^{\star} \mathcal{Y}(s, u) \right) \right),$$

and  $\{\lambda_k^{\star}: 1 \leq k \leq K\}$ ,  $\{g_k^{\star}(s): 1 \leq k \leq K\}$  are the Lagrange multipliers for the constraints (1.13b) and (1.13c), respectively, and  $g_{K+1} \equiv 0$ .

Proposition 2.2 motivates a two-step approach in which  $\lambda^*$  is obtained as the solution to a convex program that maximizes the dual function  $\varphi^*$ , and then  $g^*$  are computed through the nonlinear recursion (1.21b). Hence the larger computational challenge is computing  $\lambda^*$ . Expressions for the derivatives of  $\varphi^*$  involve means and variances of  $\mathcal{Y}(X_k)$ , which invites the application of Monte Carlo techniques when the state space is large or even uncountable—see [13] for details.

(2) Feedback. Section 3 concerns control design following steps analogous to the approach used in linear systems theory to obtain the feedback control strategy (1.20). Justification of the ACOE requires that we turn to a time-homogeneous model, meaning that  $T_k = T$  and  $\phi_k^0 = \phi^0$ , independent of k.

Even with  $r_k \equiv r$  fixed, the solution to (1.4) is not time homogeneous, but on letting  $K \to \infty$  the policies converge to a solution of an ACOE. This is equivalently expressed as the solution to a deterministic optimal control problem,

(1.22) System: 
$$\nu_{k+1} = f(\nu_k, \phi_k)$$
, Cost:  $c(\nu, \phi; r) = D_{\infty}(\hat{\nu}, \phi) + \frac{\kappa}{2} [\langle \nu, \mathcal{Y} \rangle - r]^2$ ,

where the marginals  $\{\nu_k\}$  are viewed as a state process, evolving on the simplex  $\mathcal{S}(X)$ , and  $\phi_k \in \Phi$  is regarded as an input. The system equation is of the form (1.12), but simplified because of the time-homogeneity assumptions imposed here, giving

$$f(\nu, \varphi)\big|_{x'=(s', u')} = \sum_{x \in \mathsf{X}} \nu(x) T(x, s') \varphi(u' \mid s').$$

Hence f is bilinear in the pair  $(\nu, \phi)$ . Identification and justification of the term  $D_{\infty}$  in (1.22) require further notation and analysis.

Consider the infinite-horizon objective,

(1.23) 
$$\eta^{\star}(r) = \min \limsup_{K \to \infty} \frac{1}{K} \sum_{k=1}^{K} \left[ \mathcal{D}(\nu_{k}, \nu_{k}^{0}) + \frac{\kappa}{2} \left[ \langle \nu_{k}, \mathcal{Y} \rangle - r \right]^{2} \right]$$

in which the minimum is over all  $\{\phi_k\} \subset \Phi$ . The following notational conventions are required to describe the structure of its solution:

(i) Any  $\phi \in \Phi$  defines a transition matrix  $P_{\phi}$ , and any pmf  $\pi$  that is invariant for  $P_{\phi}$  admits the decomposition

(1.24) 
$$\pi(s, u) = \phi(u \mid s)\hat{\nu}(s),$$

where  $\hat{\nu}$  is the steady-state pmf for S under this policy.

(ii) With  $\phi$  and  $\hat{\nu}$  as above, the steady-state relative entropy rate is denoted

(1.25) 
$$D_{\infty}(\hat{\nu}, \phi) := \sum_{s,u} \phi(u \mid s) \hat{\nu}(s) \log \left( \frac{\phi(u \mid s)}{\phi^{0}(u \mid s)} \right).$$

THEOREM 1.3. [infinite-horizon KLQ solution]. Suppose that the nominal transition matrix  $P^0$  has unique invariant  $pmf \pi^0$ , and fix any  $\kappa > 0$  and  $r \in \mathbb{R}$ . Then, there is a solution to (1.23) in which  $\varphi_k^* = \varphi^*$  for each k, obtained from the optimization problem

(1.26) 
$$\operatorname*{arg\,min}_{\pi,\Phi} \bigg\{ D_{\infty}(\hat{\nu}, \Phi) + \frac{\kappa}{2} [\langle \pi, \mathcal{Y} \rangle - r]^2 : \pi P_{\Phi} = \pi \bigg\}.$$

This optimization problem is convex with unique solution  $\{\pi^*, \phi^*\}$ .

The convex program (1.26) reduces to the "IPD" convex program of [37, 8, 22] as  $\kappa \to \infty$  (see discussion surrounding (1.29) in the literature review). The two convex programs are differentiated by the introduction of a quadratic cost on the marginals, so the policy  $\phi^*$  obtained from (1.26) is henceforth called the *IPD-Q* solution. Much of section 3 is devoted to obtaining approximations of this solution, as well as computational methods to obtain the exact solution.

(a) HJB solution and LQR approximation. Viewed as a deterministic optimal control problem, with system and cost given in (1.22), another solution to (1.23) is obtained as state feedback  $\phi_k^* = \mathcal{K}^*(\nu_k^*, r)$  for some mapping  $\mathcal{K}^* : \mathcal{S}(\mathsf{X}) \times \mathbb{R} \to \Phi$ . The IPD-Q solution is obtained via  $\phi^* = \mathcal{K}^*(\nu^*; r)$  with  $\nu^*$  the steady-state marginal of S under  $P_{\Phi^*}$ .

Because computation of  $\mathcal{K}^*$  is complex if |X| is large, and in anticipation of finer analysis of the performance of this policy, much of section 3.1 is devoted to "small signal" approximations.

Let  $\{x^i = (s^i, u^i) : 1 \le i \le n\}$  be an enumeration of the state space X with n = |X|. As a corollary to Propositions 3.2 and 3.3, coefficients  $\{K_{i,j}^{\star}, G_i^{\star}\}$  are constructed for which  $\phi_k^{\star}(u^i \mid s^i, r) = \phi_k(u^i \mid s^i, r) + O(r^2)$  for each i with

(1.27)

$$\phi_k(u^i \mid s^i, r) := \phi^0(u^i \mid s^i) \exp\left(\frac{1}{\phi^0(u^i \mid s^i)} \left(-\sum_j K_{i,j}^{\star} \widetilde{\nu}_k(x^j) + G_i^{\star} r\right) - \Gamma(s^i, r)\right),$$

 $\widetilde{\nu}_k(x^i) := \nu_k(x^i) - \pi^0(x^i)$  for each i, and  $\Gamma$  is a normalizing constant, defined so that  $\phi_k(\cdot \mid s^i, r)$  is a pmf on  $\mathsf{U}$  for each  $s^i, r$ .

(b) Lagrangian relaxation. A Lagrangian relaxation leads to a characterization of the IPD-Q solution in terms of a standard ACOE. Similar to (1.13b), we introduce the variable  $\gamma = \langle \nu_k, \mathcal{Y} \rangle - r$ , and let  $\lambda^* \in \mathbb{R}$  denote the Lagrange multiplier associated with this scalar constraint; it is identified in (3.4) as  $\lambda^* = \kappa[r - \langle \pi^*, \mathcal{Y} \rangle]$ .

The relative value function  $h^*$  that solves the ACOE provides a representation of the IPD-Q solution in (3.6):

$$\varphi^{\star}(u \mid s) = \varphi^{0}(u \mid s) \exp(\bar{h}^{\star}(s, u) + \lambda^{\star} \mathcal{Y}(s, u) - \Gamma^{\star}(s)),$$

with  $\bar{h}^{\star}(x) = \sum_{n'} T(x, s') h^{\star}(s')$  and  $\Gamma^{\star}(s)$  a normalizing factor.

(c) ODE solution and small signal approximation. Rather than compute  $\lambda^*$  for each r, it is argued that it is simpler to let  $\lambda$  be the independent variable. The family of relative value functions  $\{h^{\lambda}: \lambda \in \mathbb{R}\}$  solves an ordinary differential equation, whose vector field is identified in (3.10). In addition to offering a tool for exact computation, this leads to approximation of the IPD-Q solution.

These conclusions lead to several approaches to feedback control for tracking a time varying reference signal. Remember, in the following three options, the family  $\{\phi_k : k \geq 0\}$  is proposed for local decision making in a mean-field control architecture:

- 1. The feedback solution (1.17) using the collection  $\{\phi_k^{\star}(\cdot \mid \cdot; r) : k \geq 0, r \in \mathbb{R}\}$ .
- 2. The open-loop strategy  $\phi_k(\cdot | \cdot) = \phi^*(\cdot | \cdot; r_k)$ , with  $\{\phi^*(\cdot | \cdot; r) : r \in \mathbb{R}\}$  the IPD-Q solutions.
- 3. In option 2 above, it is assumed that  $r_k$  is made available to each agent, at each time k, as an external control signal. A refinement is obtained by designing a control signal  $\{\zeta_k : k \geq 0\}$  based on filtering measurements, such as error feedback,

(1.28) 
$$\zeta_k = \sum_{i=0}^k g_{k-i} e_i, \quad k \ge 0, \qquad e_i = r_i - \langle \nu_i, \mathcal{Y} \rangle.$$

The randomized decision rule for each agent is then  $\phi_k(\cdot | \cdot) = \phi^*(\cdot | \cdot; \zeta_k)$  with  $\{\phi^*(\cdot | \cdot; \zeta) : \zeta \in \mathbb{R}\}$  the IPD-Q solutions. The linearized dynamics described in Proposition 3.5 can aid in the design of the filter in (1.28).

(3) Application to demand dispatch. The original motivation for the research surveyed here is application to distributed control of power systems. The term demand dispatch was introduced in the conceptual article [7] to describe the possibility of distributed intelligence in electric loads, designed so that the population would help provide supply-demand balance in the power grid.

The numerical results surveyed in section 4 illustrate the application of KLQ to control a large population of residential loads. As expected, tracking error can be made arbitrarily small with large  $\kappa > 0$ , provided the reference signal is feasible.

It is found in numerical experiments that the histograms defining the state of the mean-field model rapidly "forget" their initial conditions—see the full arXiv version [13] for details.

#### 1.6. Literature review.

Mean-field control. The optimization problem (1.3) is inspired by mean-field game theory [31, 28, 29, 10, 26] (see [16, 17, 10, 42] for recent surveys).

Mean-field control differs from mean-field game theory only because of greater control at the microscopic layer: we do not assume that an individual in the population is free to optimize based on its local objective function, so we avoid the fragility of Nash equilibria. This description is similar to *ensemble control* in physics (see [32] for history), and many in the power systems area opt for this term rather than mean-field control (see [23, 22] and their references).

Demand dispatch. The goal of demand dispatch is to modify the behavior of loads so that their aggregate power consumption tracks a reference signal  $\{r_k\}$  that is synthesized by a balancing authority (BA). Randomized control techniques have been proposed in [35, 43, 37, 1, 23, 4] based on various control architectures.

The following control strategy is common to the approaches described in [37, 22]. It is assumed that a family of transition matrices  $\{P_{\zeta}: \zeta \in \mathbb{R}\}$  is available at each load. A sequence  $\{\zeta_0, \zeta_1, \ldots\}$  is broadcast from the BA, based on measurements of the grid, and at time k the ith load transitions according to this law:

$$P\{X_{k+1}^i = x' \mid X_k^i = x, \zeta_k = \zeta\} = P_{\zeta}(x, x').$$

The feedback solution (1.28) was proposed in [37] and tested in this and later research using  $e_i = r_i - \langle \nu_i^{\mathcal{N}}, \mathcal{Y} \rangle$  [22].

*IPD.* The paper [37] reinterprets the control solution of [44] as a technique to create the family  $\{P_{\zeta}\}$  through the solution to the nonlinear program:

(1.29) 
$$P_{\zeta} := \arg \max \left\{ \zeta \langle \pi, \mathcal{Y} \rangle - \mathcal{R}(P \| P^{0}) \right\}, \qquad \zeta \in \mathbb{R},$$

where  $\mathcal{R}$  denotes the rate function of Donsker and Varadhan [25, 30],

(1.30) 
$$\mathcal{R}(P||P^0) := \sum_{x,x'} \pi(x) P(x,x') \log \left( \frac{P(x,x')}{P^0(x,x')} \right),$$

in which  $\pi$  is the invariant pmf for P. The maximum in (1.29) is over all  $(\pi, P)$  subject to the invariance constraint  $\pi P = \pi$  [37, 8]. The convex program (1.29) is called the individual perspective design (IPD) in [8].

Hence IPD-Q may be interpreted as a new approach to designing  $\{P_{\zeta}\}$ .

The finite-horizon version of (1.29) is also considered in [37, 8], similar to the KLQ formulation:

$$(1.31) p^{\zeta} := \underset{p}{\operatorname{arg\,max}} \left\{ \zeta \mathsf{E}_{p} \left[ \sum_{k=1}^{K} \mathcal{Y}(x_{k}) \right] - D(p \| p^{0}) \right\}.$$

Provided the entries of  $T_k(x,s)$  take on only binary values, the finite-horizon IPD solution is obtained as a tilting of the nominal model:

$$(1.32) \ p^{\zeta}(\vec{x}) = p^{0}(\vec{x}) \exp\left(\zeta \sum_{k=1}^{K} \mathcal{Y}(x_{k}) - \Lambda(\zeta)\right) \quad \text{with } \Lambda(\zeta) \ a \ normalizing \ constant.$$

KLQ and optimal transport. Extensions of the KLQ objective will likely provide useful relaxations of the classical optimal transport problem, in which the goal is to steer  $p^0$  to a given target pmf  $p^*$  [39, 21]. Rather than match the target pmf, we might match M generalized moments, minimizing  $D(p||p^0)$  subject to  $\langle p, \mathcal{G}_i \rangle = \langle p^*, \mathcal{G}_i \rangle$  for each i, with  $\mathcal{G}_i : \mathsf{X}^{K+1} \to \mathbb{R}$ .

A special case is the tracking problem,

(1.33) 
$$\min_{p} \left\{ D(p \| p^0) \text{ subject to } \mathsf{E}_p \big[ \mathcal{Y}(X_k) \big] = r_k \,, \quad 1 \le k \le K \right\}.$$

This optimization problem is proposed in [23, section 5], along with the explicit solution

(1.34) 
$$p^{\star}(\vec{x}) = p^{0}(\vec{x}) \exp\left(\sum_{k=1}^{K} \beta_{k} \mathcal{Y}(x_{k}) - \Lambda(\beta)\right)$$

in which  $\beta \in \mathbb{R}^K$  are Lagrange multipliers corresponding to the K constraints, and  $\Lambda(\beta)$  a normalizing constant.

The convex program formulation (1.13) has many advantages. First, (1.13) is always feasible, while feasibility of (1.33) requires conditions on  $p^0$  and  $\{r_k\}$ . Theorem 1.2 requires no assumptions on the model or reference signal. Flexibility in choice of  $\kappa$  allows for *learning* the characteristics of an "expensive" reference signal. It is anticipated that the penalty parameter  $\kappa$  can be used to make trade-offs between tracking performance and robustness to modeling error: robustness and sensitivity analysis will be a topic of future research.

Finally, as assumed to obtain the representation (1.32), the formula (1.34) is meaningful only when  $T_k(x,s)$  take on only binary values. A goal of the research surveyed in this paper is to remove this restriction.

The similarity between (1.32) and (1.34) is not accidental but follows from an alternative interpretation of the IPD design (1.31). For a scalar  $r_0 \in \mathbb{R}$ , consider the constrained optimization problem

$$(1.35) \qquad \max_{p} \left\{ -D(p\|p^0) \right\} \quad \text{subject to } \mathsf{E}_{p} \left[ \sum_{k=1}^{K} \mathcal{Y}(x_k) \right] = K r_0 \,, \quad 1 \leq k \leq K.$$

The dual function  $\varphi^* \colon \mathbb{R} \to \mathbb{R}$  is defined by

$$\varphi^{\star}(\lambda) = \max_{p} \left\{ \lambda \mathsf{E}_{p} \left[ \sum_{k=1}^{K} \mathcal{Y}(x_{k}) \right] - D(p \| p^{0}) \right\} - \lambda K r_{0},$$

where  $\lambda \in \mathbb{R}$  is a Lagrange multiplier. It is evident that the optimizer  $p^{*\lambda}$  is an IPD solution for each  $\lambda$ . Consequently, for each  $\zeta$ , the IPD solution (1.31) also solves (1.35) for some scalar  $r_0(\zeta)$ .

Contributions. Most of the contributions were surveyed in section 1.5. The main contribution of this paper is the discovery of hidden convexity in the nonlinear program (1.13), which leads to structure for the optimal solution in Theorem 1.2. Properties of the dual surveyed in Theorem 2.1 lead to computational techniques for this new class of optimal control formulations; see Proposition 2.2 and its corollary. The application of these techniques to the infinite-horizon setting in section 3 is novel, and the main results surveyed there are new.

Portions of the results reported here were summarized in the conference article [15]. In this preliminary work, the transition matrix  $T_k$  was assumed *deterministic*, so that all randomness arose from the randomized policy. All of the results in this paper allow for general Markovian dynamics.

Extensions to resource allocation are summarized in [14]. More on these topics may be found in the first author's Ph.D. dissertation [12].

**Organization.** The remainder of this paper is organized as follows. Section 2 describes a relaxation technique motivated by the desire to reduce computational complexity, along with a full analysis of the convex program (1.13) and its dual. Section 3 contains extensions to the infinite-horizon setting. Results from numerical experiments are collected together in section 4. Conclusions and directions for future research are contained in section 5.

## 2. Kullback-Leibler-quadratic optimal control.

2.1. Subspace relaxation. A relaxation of the convex program (1.13) is described here. Motivation is most clear from consideration of distributed control of a collection of residential water heaters. These loads are valuable as sources of virtual energy storage since they in fact are energy storage devices (in the form of heat rather than electricity) and are also highly flexible. Flexibility comes in part from their extremely nonsymmetric behavior: a typical unit may be on for just five minutes and off continuously for more than six hours. The intersampling time at the load should be far less than five minutes to obtain a reliable model for control.

On the other hand, it is valuable for the time horizon to be on the order of several hours. For example, peak-shaving is more effective when water heaters have advance warning to preheat the water tanks. To obtain a useful control solution will thus require a very large value of K in (1.13). To reduce complexity, an approach is proposed here based on lossy compression of  $\{r_k\}$  using transform techniques.

The transformations are based on a collection of functions  $\{w_n: 1 \leq n \leq N\}$ , with  $w_n: \{0,1,\ldots,K\} \to \mathbb{R}$  for each n, and  $N \ll K$ . The transformed signal is the N-dimensional vector  $\hat{r}$  with  $\hat{r}_n = \sum_k w_n(k) r_k$  for each n, and the transformed function on  $\mathsf{X}^{K+1}$  is denoted  $\widehat{\mathcal{Y}}_n(\vec{x}) = \sum_{k=1}^K w_n(k) \mathcal{Y}(x_k)$  for  $1 \leq n \leq N$ .

The goal is to achieve the approximation  $\langle p, \widehat{\mathcal{Y}}_n \rangle \approx \hat{r}_n$  for each n, while maintaining  $p \approx p^0$ . For example, a Fourier series can be used, with frequency  $\omega > 0$ , and N is necessarily odd:  $w_n(k) \in \{1, \sin(\omega m k), \cos(\omega m k) : 1 \leq m \leq (N-1)/2\}$ . The degenerate family is defined using N = K and

(2.1) 
$$w_n(k) = \mathbb{I}\{n = k\}, \quad 1 \le n, k \le K.$$

The optimal control problem with subspace relaxation is defined as

(2.2a) 
$$J^{\star}(\nu_0^0) := \min_{\nu,\gamma} \sum_{k=1}^K \mathcal{D}(\nu_k, \nu_k^0) + \frac{\kappa}{2} \sum_{n=1}^N \gamma_n^2$$

(2.2b) s.t. 
$$\gamma_n = \langle p, \widehat{\mathcal{Y}}_n \rangle - \hat{r}_n, \quad 1 \le n \le N,$$

(2.2c) 
$$\sum_{u'} \nu_k(s', u') = \sum_{s, u} \nu_{k-1}(s, u) T_{k-1}(x, s'), \qquad 1 \le k \le K, \ s' \in S.$$

This reduces to (1.13) in the degenerate case (2.1).

The theory that follows is based in part on a relaxation of the dynamical constraints (2.2c), through the introduction of a Lagrange multiplier for each k. This is precisely the first step in the construction of the Hamiltonian in the minimum principle approach to optimal control [33].

**2.2. Duality.** Structure for the solution of (2.2) will be obtained by consideration of a dual, in which  $\lambda \in \mathbb{R}^N$  and  $g \in \mathbb{R}^{K \times J^*}$  denote the vectors of Lagrange multipliers for the first and second sets of constraints, respectively. The matrix g is interpreted as a sequence of functions  $g_k \colon S \to \mathbb{R}$  that are entirely analogous to the costate variables in the minimum principle (the Lagrange multipliers for the dynamical constraints) [33].

The Lagrangian is denoted

$$\mathcal{L}(\nu, \gamma, \lambda, g) = \sum_{k=1}^{K} \mathcal{D}(\nu_{k}, \nu_{k}^{0}) + \frac{\kappa}{2} \sum_{n=1}^{N} \gamma_{n}^{2} + \sum_{n=1}^{N} \lambda_{n} \left( \gamma_{n} + \sum_{k=1}^{K} w_{n}(k) \left[ r_{k} - \langle \nu_{k}, \mathcal{Y} \rangle \right] \right)$$

$$+ \sum_{k=1}^{K} \sum_{s'} \left( \sum_{u'} \nu_{k}(s', u') - \sum_{s, u} \nu_{k-1}(s, u) T_{k-1}(x, s') \right) g_{k}(s')$$

$$(2.3)$$

and the dual function is defined to be its minimum,  $\varphi^*(\lambda, q) := \min_{\nu, \gamma} \mathcal{L}(\nu, \gamma, \lambda, q)$ .

The dual of the optimization problem (2.2) is defined as the maximum of the dual function  $\varphi^*$  over  $\lambda$  and g (see [33] for a complete and accessible treatment of this theory). We will see that there is no duality gap, so that for a quadruple  $(\nu^*, \gamma^*, \lambda^*, g^*)$ ,

$$J^{\star}(\nu_0^0) = \mathcal{L}(\nu^{\star}, \gamma^{\star}, \lambda^{\star}, g^{\star}) = \varphi^{\star}(\lambda^{\star}, g^{\star}).$$

In the following subsections a representation of the dual function is obtained that is suitable for optimization, which results in a valuable representation for the optimal policy. Properties of the dual function are contained in Theorem 2.1 and Proposition 2.2 that follow. The statement of these results requires additional notation: define a function  $\mathcal{T}_k^{\lambda} \colon \mathbb{R}^{|\mathsf{S}|} \to \mathbb{R}^{|\mathsf{S}|}$ , for  $f \colon \mathsf{S} \to \mathbb{R}$  and  $\lambda \in \mathbb{R}^N$ , via

$$\mathcal{T}_k^{\lambda}(f;s) = \log \left( \sum_{u} \Phi_k^0(u \mid s) \exp \left( \sum_{s'} T_k(x,s') f(s') + \check{\lambda}_k \mathcal{Y}(s,u) \right) \right), \quad s \in \mathsf{S},$$

where  $\check{\lambda}_k = \sum_{n=1}^N \lambda_n w_n(k)$ . The maximum of the dual function over g is denoted

$$\varphi^*(\lambda) := \max_g \varphi^*(\lambda, g) = \varphi^*(\lambda, g^{\lambda}),$$

where  $g^{\lambda}$  is a maximizer,  $g^{\lambda} \in \arg \max_g \Phi^{\star}(\lambda, g)$ . It is shown in Proposition 2.2 that the vector valued function  $g^{\lambda}$  satisfies the recursion

$$(2.5) g_k^{\lambda} = \mathcal{T}_k^{\lambda}(g_{k+1}^{\lambda}), \quad 1 \le k \le K, \quad \text{where} \quad g_{K+1}^{\lambda} \equiv 0.$$

This forms part of the proof of Theorem 2.1, with complete details postponed to Appendix B.

THEOREM 2.1. There exists a maximizer  $\{\lambda_n^{\star}, g_k^{\star} : 1 \leq n \leq N, 1 \leq k \leq K\}$  for  $\varphi^{\star}$ , and there is no duality gap:  $\varphi^{\star}(\lambda^{\star}, g^{\star}) = J^{\star}(\nu_0^0)$ . The optimal policy is obtained from  $\{g_k^{\star}\}$  via

(2.6) 
$$\phi_k^{\star}(u \mid s) = \phi_k^0(u \mid s) \exp\left(\sum_{s'} T_k(x, s') g_{k+1}^{\star}(s') + \check{\lambda}_k^{\star} \mathcal{Y}(s, u) - g_k^{\star}(s)\right),$$

$$where \ g_k^{\star}(s) = \mathcal{T}_k^{\lambda}(g_{k+1}^{\star}; s) \ for \ 1 \le k \le K, \ and \ g_{K+1}^{\star} \equiv 0,$$

and  $\{\check{\lambda}_k^{\star}\}$  are obtained from  $\{\lambda_n^{\star}\}$  via (2.4).

The proof of the following is also contained in Appendix B. Denote for each k,

(2.7) 
$$G_k^{\lambda}(x) = \sum_{s} T_{k-1}(x, s) g_k^{\lambda}(s).$$

PROPOSITION 2.2. The following hold for the dual of (2.2): for each  $\lambda \in \mathbb{R}^N$ ,

- (i) a maximizer  $g^{\lambda}$  is given by (2.5);
- (ii) the maximum of the dual function over g is the concave function

(2.8) 
$$\varphi^{\star}(\lambda) = \lambda^T \hat{r} - \frac{1}{2\kappa} ||\lambda||^2 - \langle \nu_0^0, G_1^{\lambda} \rangle;$$

(iii) the function (2.8) is continuously differentiable, and

(2.9) 
$$\frac{\partial}{\partial \lambda_n} \varphi^*(\lambda) = \hat{r}_n - \frac{1}{\kappa} \lambda_n - \sum_{k=1}^K w_n(k) \langle \nu_k^{\lambda}, \mathcal{Y} \rangle, \qquad 1 \le n \le N,$$

where  $\{\nu_k^{\lambda}\}$  is the sequence of marginals obtained from the randomized policy defined in (2.6), substituting  $\{g_k^{\star}\}$  by  $\{g_k^{\lambda}\}$  defined in (i).

To conclude this section, we provide representations of the log-likelihood ratio,  $L(\vec{x})$ , relative entropy  $D(p^{\lambda}||p^{0})$ , and primal objective function for the pmf  $p^{\lambda} \in \mathcal{S}(\mathsf{X}^{K+1})$  obtained from the randomized policy defined in (2.6), substituting  $\{g_{k}^{\star}\}$  by  $\{g_{k}^{\lambda}\}$  defined in Proposition 2.2, part (i). We defer to [13] for the proof of the following.

Corollary 2.3. The following hold for all  $\{\check{\lambda}_k, g_k^{\lambda} : 1 \leq k \leq K\}$ :

(i) The log-likelihood ratio can be expressed as

(2.10) 
$$L(\vec{x}) = \sum_{k=1}^{K} \{ \Delta_k(x_{k-1}, s_k) + \check{\lambda}_k \mathcal{Y}(x_k) \} - G_1^{\lambda}(x_0),$$

where for each k (recalling  $x_k = (s_k, u_k)$ ),

(2.11) 
$$\Delta_k(x_{k-1}, s_k) = G_k^{\lambda}(x_{k-1}) - g_k^{\lambda}(s_k).$$

(ii) The relative entropy is given by

(2.12) 
$$D(p^{\lambda}||p^{0}) = \sum_{k=1}^{K} \check{\lambda}_{k} \langle \nu_{k}^{\lambda}, \mathcal{Y} \rangle - \langle \nu_{0}^{0}, G_{1}^{\lambda} \rangle.$$

(iii) The value of the primal is given by

(2.13a) 
$$J(p^{\lambda}, \nu_0^0) := D(p^{\lambda} || p^0) + \frac{\kappa}{2} \sum_{n=1}^N \gamma_n^2$$

$$= -\langle \nu_0^0, G_1^{\lambda} \rangle + \sum_{k=1}^K \check{\lambda}_k \langle \nu_k^{\lambda}, \mathcal{Y} \rangle + \frac{\kappa}{2} \sum_{n=1}^N \gamma_n^2$$

with 
$$\gamma_n = \langle p^{\lambda}, \widehat{\mathcal{Y}}_n \rangle - \hat{r}_n$$
.

The stochastic process  $\{\Delta_k(X_{k-1}, S_k)\}$  is a martingale difference sequence; it vanishes when nature is deterministic, reducing to the solution obtained in [15].

**3. Feedback formulations.** We now turn to the IPD-Q convex program (1.26). It is assumed throughout this section that  $T_k = T$  and  $\phi_k^0 = \phi^0$ , independent of k.

The relationship between IPD-Q and (1.23) will be clear after justification of the term  $D_{\infty}(\hat{\nu}, \phi)$  defined in (1.25). Consider any  $\phi \in \Phi$ , which gives rise to a Markov chain with transition matrix  $P_{\phi}$ . The relative entropy (1.15) was previously expressed as a sum over  $\vec{x} \in X^{K+1}$  in (1.15). The notation  $D(p||p^0) = D^K(p||p^0)$  is required in the following, since K is a variable in (1.23).

PROPOSITION 3.1. Suppose that p is obtained using the policy  $\phi$ , and initial pmf  $\nu_0^0$  common with  $p^0$ . Suppose moreover that  $P_{\phi}$  has a unique invariant pmf  $\pi$ . Then,

$$D_{\infty}(\hat{\nu}, \varphi) = \lim_{K \to \infty} \frac{1}{K} \sum_{k=1}^{K} \mathcal{D}(\nu_k, \nu_k^0) = \lim_{K \to \infty} \frac{1}{K} D^K(p \| p^0) = \mathcal{R}(P_{\varphi} \| P^0),$$

where  $\mathcal{R}$  denotes the rate function (1.30) using  $P = P_{\Phi}$ :

$$\mathcal{R}(P_{\Phi} \| P^0) = \sum_{x.x'} \pi(x) P_{\Phi}(x, x') \log \left( \frac{P_{\Phi}(x, x')}{P^0(x, x')} \right).$$

*Proof.* The proof of the first identity begins with

$$\frac{1}{K} \sum_{k=1}^{K} \mathcal{D}(\nu_{k}, \nu_{k}^{0}) = \frac{1}{K} \sum_{k=1}^{K} \mathsf{E}[F(X_{k}))]$$

with  $F(x) = \log[\phi_k(u \mid s)/\phi_k^0(u \mid s)]$  for  $x = (s, u) \in X$ . The average converges to  $D_{\infty}(\hat{\nu}, \phi)$  as  $K \to \infty$  since the invariant pmf  $\pi$  is unique.

The distinct approaches to optimal control pursued in this section follow the distinct approaches to optimal control in general, via the HJB equations and optimal control via the minimum principle:

(i) In section 3.1 IPD-Q is interpreted as a solution to an HJB equation, which results in a solution in state feedback form,  $\phi_k^* = \mathcal{K}^*(\nu_k^*, r)$ , for some mapping  $\mathcal{K}^* \colon \mathcal{S}(\mathsf{X}) \times \mathbb{R} \to \Phi$ . The solution to IPD-Q is  $\phi^*(u \mid s) = \mathcal{K}^*(\nu^r, r)$ , in which  $\nu^r$  is the steady-state marginal for S under the IPD-Q policy. Computation

of  $\mathcal{K}^*$  may be difficult if the state space is large. An LQR approximation is proposed, justified for small |r|, and the approximation (1.27) may also be found at the close of section 3.1.

- (ii) The approach taken in section 3.2 is in essence the infinite-K limit of the approach taken in section 2.2 which, as noted following (2.2), is the minimum principle approach. It is well known that this approach provides only an open-loop solution.
- **3.1. HJB approach.** The solution to the optimal control problem (1.22) may be characterized using techniques from deterministic optimal control theory.

The ACOE holds for deterministic systems, precisely as reviewed in section 1.4 for the linear quadratic problem:

(3.1) 
$$\eta^{\star} + \mathcal{H}^{\star}(\nu) = \min_{\Phi} \left\{ c(\nu, \Phi; r) + \mathcal{H}^{\star}(f(\nu, \Phi)) \right\}, \qquad \nu \in \mathcal{S}(\mathsf{X}).$$

with  $c(\nu, \phi; r)$  defined in (1.22),  $\mathcal{H}^* \colon \mathcal{S}(\mathsf{X}) \to \mathbb{R}$  the relative value function, and  $\eta^*$  the optimal average cost. The minimizer  $\phi^*$  defines  $\phi_k^* = \mathcal{K}^*(\nu_k^*, r)$ .

We are not aware of solution techniques for this instance of the ACOE, beyond the standard value iteration algorithm or other generic approaches.

The relative value function  $\mathcal{H}^*$  and feedback law  $\mathcal{K}^*$  can be approximated through a small signal linearization of the dynamics, and a quadratic approximation of the cost. We begin with an approximation for the latter.

The proof of Proposition 3.2 follows from the definition (1.25) and a Taylor's series approximation of the logarithm. For any  $\phi$ , denote by  $\widetilde{\phi}(u \mid s) := \phi(u \mid s) - \phi^0(u \mid s)$  the deviation.

PROPOSITION 3.2. Suppose that  $\langle \pi^0, \mathcal{Y} \rangle = 0$ . Then, the cost function (1.22) is nearly quadratic in deviations,

$$c(\nu, \mathbf{\varphi}; r) = \|\widetilde{\mathbf{\varphi}}\|_R^2 + \frac{\kappa}{2} (y - r)^2 + O(\|\widetilde{\mathbf{\varphi}}\|_R^3)$$

in which  $y = \sum_x [\nu(x) - \pi^0(x)] \mathcal{Y}(x)$ , and  $\|\widetilde{\varphi}\|_R^2 = \sum_{s,u} \frac{\hat{\nu}(s)}{\Phi^0(u|s)} \widetilde{\varphi}(u \mid s)^2$ .

Approximation of the dynamics by a linear system is justified when |r| is small, and  $\nu_0^0 \approx \pi^0$ , the invariant pmf for  $P^0$ . The corresponding stationary pmf for S is denoted  $\hat{\nu}^0$  (recall (1.24)).

Let  $X = \{x^i : 1 \le i \le n\}$  with n = |X|. The LQR approximation has state denoted  $\widetilde{\mathcal{X}}_k$  and input  $\widetilde{\mathcal{U}}_k$  at time k, with  $\widetilde{\mathcal{X}}_k^i$  an approximation of  $\widetilde{\nu}_k(x^i) := \nu_k(x^i) - \pi^0(x^i)$ , and  $\widetilde{\mathcal{U}}_k^i$  an approximation of  $\widetilde{\phi}_k^i(u^i \mid s^i) := \phi_k^i(u^i \mid s^i) - \phi^0(u^i \mid s^i)$ . The definition of the linearization is a system model of the form (1.18),

$$\widetilde{\mathcal{X}}_{k+1} = A\widetilde{\mathcal{X}}_k + B\widetilde{\mathcal{U}}_k, \qquad \widetilde{\mathcal{Y}}_k = C^T\widetilde{\mathcal{X}}_k,$$

in which  $\widetilde{\mathcal{Y}}_k$  is an approximation of  $\langle \widetilde{\nu}_k, \mathcal{Y} \rangle$ . Expressions for the  $n \times n$  matrices A and B, and the n-dimensional column vector C, are provided in the following.

Proposition 3.3. The small signal approximation holds with

$$A_{i,j} = P^0(x^j, x^i), \quad B_{i,j} = \mathbb{I}\{i = j\}\hat{\nu}^0(s^j), \quad C_i = \mathcal{Y}(x^i), \quad 1 \le i, j \le n.$$

*Proof.* The expression for C is by definition of  $\widetilde{\mathcal{Y}}_k$ . The other matrices are obtained through the standard first-order Taylor series approximations:

$$A_{i,j} := \frac{\partial}{\partial \nu^j} f_i(\nu, \varphi) \Big|_{\nu = \pi^0, \varphi = \varphi^0} = P^0(x^j, x^i)$$

with  $\nu^j = \nu(x^j)$ .

The input  $\mathcal{U}_k$  is an *n*-dimensional column vector, so that *B* is an  $n \times n$  matrix. It is obtained from the Taylor series approximation,

$$B_{i,j} := \frac{\partial}{\partial \varphi^j} f_i(\nu, \varphi) \Big|_{\nu = \pi^0, \varphi = \varphi^0} = \mathbb{I}\{i = j\} \sum_{x \in \mathsf{X}} \pi^0(x) T(x, s^j),$$

where  $\phi^j = \phi(u^j \mid s^j)$ . By invariance of  $\pi^0$  it follows that B is diagonal, with ith diagonal entry equal to  $\hat{\nu}^0(s^i)$ .

Propositions 3.2 and 3.3 imply that for small r, the solution to the nonlinear optimal control problem is approximated by the average-cost LQR solution using

$$c(\widetilde{\mathcal{X}},\widetilde{\mathcal{U}};r) = \|\widetilde{\mathcal{U}}\|_R^2 + \frac{\kappa}{2}(\widetilde{\mathcal{Y}}-r)^2\,, \quad \widetilde{\mathcal{Y}} = C^T\widetilde{\mathcal{X}},$$

giving  $\widetilde{\mathcal{U}}_k = -K^*\widetilde{\mathcal{X}}_k + G^*r$ , with gain matrices  $K^*$   $(n \times n)$  and  $G^*$   $(n \times 1)$ . This leads to the policy approximation. Write  $\widetilde{\mathcal{U}}_k = \mathcal{U}_k - \mathcal{U}_k^0$  with  $\mathcal{U}_k^0$  the vector representation of the nominal policy. The ith entry of the input is expressed

$$\begin{split} \mathcal{U}_k^i &= \boldsymbol{\Phi}^0(\boldsymbol{u}^i \mid \boldsymbol{s}^i) + [-K^\star \widetilde{\mathcal{X}}_k + G^\star \boldsymbol{r}]_i \\ &= \boldsymbol{\Phi}^0(\boldsymbol{u}^i \mid \boldsymbol{s}^i) \left[ 1 + \frac{1}{\boldsymbol{\Phi}^0(\boldsymbol{u}^i \mid \boldsymbol{s}^i)} \left( -[K^\star \widetilde{\mathcal{X}}_k]_i + G_i^\star \boldsymbol{r} \right) \right]. \end{split}$$

This implies the small signal approximation (1.27). It is conjectured that (1.27) is within  $O(r^2)$  of optimal (in terms of the objective in (1.26)).

- **3.2.** Minimum principle approach. As previously observed, the optimization problem (1.26) falls outside of traditional MDP theory:
  - (i) The control cost is absent and is replaced by a cost on the randomized policy.
  - (ii) A quadratic cost on  $\pi$  appears, rather than linear as anticipated in the LP formulations of MDPs.

An MDP model is constructed here through a series of steps, with the first step addressing (i). For this it is natural to view the input as an element of the simplex  $\mathcal{S}(\mathsf{U})$ . This is not the same setting as section 3.1: in this subsection, the notation  $\phi(\cdot \mid s)$  is interpreted as static state feedback from state s to input  $\phi(\cdot \mid s) \in \mathcal{S}(\mathsf{U})$ .

To remove the quadratic cost on  $\pi$  requires a Lagrangian relaxation, similar to what was used in section 2. For  $\lambda \in \mathbb{R}$  denote

$$(3.2) \qquad [\pi^{\lambda}, \phi^{\lambda}, \gamma^{\lambda}] = \operatorname*{arg\,min}_{\pi, \phi, \gamma} \left\{ \mathcal{D}(\phi) + \frac{\kappa}{2} \gamma^2 + \lambda [\gamma - \langle \pi, \mathcal{Y} \rangle + r] : \pi P_{\phi} = \pi \right\}.$$

For each  $\lambda$  this is viewed as a standard average-cost optimal control problem with state process S. The controlled transition matrix and cost function are defined by

$$\begin{split} T_{\mu}(s,s') &:= \sum_{u} \mu(u) T((u,s),s') \,, & s,s' \in \mathsf{S} \,, \quad \mu \in \mathcal{S}(\mathsf{U}) \,, \\ c(s,\mu) &:= \sum_{u} \mu(u) \left[ \log \left( \frac{\mu(u)}{\Phi^0(u \, | \, s)} \right) - \lambda \mathcal{Y}(s,u) \right], \quad s \in \mathsf{S} \,, \quad \mu \in \mathcal{S}(\mathsf{U}) \,. \end{split}$$

Under any policy  $\phi \in \Phi$  the resulting process **S** is Markovian. With a slight abuse of notation, its transition matrix is denoted

$$T_{\Phi}(s,s') := \sum_{u} \Phi(u \mid s) T((u,s),s')$$

and the cost as a function of s under this policy is denoted

$$c_{\Phi}(s) = \sum_{u} \Phi(u \mid s) c(s, \Phi(u \mid s)) = \sum_{u} \Phi(u \mid s) \left[ \log \left( \frac{\Phi(u \mid s)}{\Phi^{0}(u \mid s)} \right) - \lambda \mathcal{Y}(s, u) \right], \quad s \in S$$

The solution to (3.2) gives  $\gamma^{\lambda} = -\lambda/\kappa$  and

$$[\pi^{\lambda}, \varphi^{\lambda}] = \operatorname*{arg\,min}_{\hat{\nu}, \varphi} \left\{ \sum_{s} \hat{\nu}(s) c_{\varphi}(s) : \hat{\nu} T_{\varphi} = \hat{\nu} \right\}.$$

This is a standard MDP formulation, in which the optimization over feedback laws  $\phi$  is explicit.

The Lagrange multiplier  $\lambda$  is treated as the independent parameter rather than r. This is justified through the correspondence  $\gamma^{\lambda} = -\lambda/\kappa$ , and the following definition imposes complementary slackness:

(3.4) 
$$r^{\lambda} = -\gamma + \langle \pi^{\lambda}, \mathcal{Y} \rangle = \langle \pi^{\lambda}, \mathcal{Y} \rangle + \lambda/\kappa.$$

As  $\lambda$  ranges from  $-\infty$  to  $+\infty$ , so do the values of  $r^{\lambda}$  because  $\langle \pi^{\lambda}, \mathcal{Y} \rangle$  is bounded and continuous in  $\lambda$ .

Continuity of  $\langle \pi^{\lambda}, \mathcal{Y} \rangle$  and other conclusions are obtained from prior research (in particular [9]), because the optimization problem (3.3) is identical to the IPD optimization problem (1.29), in which  $\zeta$  is replaced by  $\lambda$ .

To match the setting of [9], denote the one-step reward as the negative of cost,  $\varrho(s, \phi) = -c(s, \phi)$ . Based on the foregoing, the solution to (3.3) is characterized by the average reward optimality equation

(3.5) 
$$\xi^{\lambda} + h^{\lambda}(s) = \max_{\Phi} \left\{ \varrho(s, \Phi) + \sum_{s'} T_{\Phi}(s, s') h^{\lambda}(s') \right\}.$$

The maximizer provides a representation for the optimal policy similar to (1.21a):

(3.6) 
$$\varphi^{\lambda}(u \mid s) = \varphi^{0}(u \mid s) \exp(\bar{h}^{\lambda}(s, u) + \lambda \mathcal{Y}(s, u) - \Gamma^{\lambda}(s)),$$

with  $\bar{h}^{\lambda}(x) = \sum_{u'} T(x, s') h^{\lambda}(s')$  and  $\Gamma^{\lambda}(s)$  the normalizing factor,

(3.7) 
$$\Gamma^{\lambda}(s) = \log \sum_{u} \Phi^{0}(u \mid s) \exp(\bar{h}^{\lambda}(s, u) + \lambda \mathcal{Y}(s, u)).$$

**ODE solution.** The reader is referred to [9] for full details on this solution technique to compute the solution to (3.5). The main ideas are recalled here, in part because they are required in a small signal approximation.

It is shown in this prior work that the relative value functions can be constructed so that they are continuously differentiable in  $\lambda$ . Letting  $H^{\lambda} = \frac{d}{d\lambda}h^{\lambda}$ , we obtain

$$\overline{\mathcal{Y}}^{\lambda} + H^{\lambda}(s) = \mathcal{Y}_s + \sum_{s'} T_{\lambda}(s, s') H^{\lambda}(s') \,, \quad s \in \mathbb{S} \,, \ \lambda \in \mathbb{R} \,,$$

in which  $T_{\lambda} = T_{\phi^{\lambda}}$ ,

(3.8) 
$$\mathcal{Y}_s := \frac{d}{d\lambda} \varrho(s, \phi) = \sum_{u} \phi(u \mid s) \mathcal{Y}(s, u) \quad \text{and} \quad \overline{\mathcal{Y}}^{\lambda} = \frac{d}{d\lambda} \xi^{\lambda}.$$

This fixed point equation is known as Poisson's equation, whose solution is often expressed  $H^{\lambda} = Z_{\lambda} \mathcal{Y}$  with  $Z_{\lambda}$  known as the fundamental matrix (obtained as a simple matrix inverse). Also obtained is

(3.9) 
$$\overline{\mathcal{Y}}^{\lambda} = \sum_{x} \pi^{\lambda}(x) \mathcal{Y}(x),$$

where  $\pi^{\lambda}(s,u) = \hat{\nu}^{\lambda}(s) \Phi^{\lambda}(u \mid s)$ , with  $\hat{\nu}^{\lambda}$  the unique invariant pmf for  $T_{\lambda}$ .

This defines the ODE solution for the family of relative value functions

$$(3.10) \qquad \frac{d}{d\lambda}h_{\lambda} = Z_{\lambda}\mathcal{Y}$$

with boundary condition  $h^{\lambda} \equiv 0$  when  $\lambda = 0$ . The right-hand side depends on  $h_{\lambda}$ through  $Z_{\lambda}$ , but the dependency is smooth.

Small signal approximation. The small signal approximation here is defined in a setting similar to Proposition 3.3: it is assumed that the reference signal is small in magnitude, and that  $r \equiv 0$  achieves zero cost in (1.26). This holds if  $\sum \pi^0(x)\mathcal{Y}(x) = 0$ , which will be assumed henceforth.

A slight change in notation is required here, as compared to section 3.1:  $X_k$  and  $U_k$  are n-dimensional column vectors that denote the exact deviation,  $X_k^i := \widetilde{\nu}_k(x^i)$ and  $\widetilde{U}_k^i = \widetilde{\Phi}_k(u^i \mid s^i)$  for each i and k. The approximation requires the following

- (i)  $\zeta_{\lambda}^{2} = \frac{d^{2}}{d\lambda^{2}} \xi^{\lambda}$  for  $\lambda \in \mathbb{R}$ . (ii)  $\bar{H}^{\lambda} = \frac{d}{d\lambda} \bar{h}^{\lambda}$ ,  $\bar{H}_{s}^{\lambda} = \sum_{u} \phi^{0}(u \mid s) \bar{H}^{\lambda}(s, u)$ ,  $\mathcal{Y}_{s} = \sum_{u} \phi^{0}(u \mid s) \mathcal{Y}(s, u)$ .
- (iii)  $\Lambda(x) = \bar{H}^0(x) + \mathcal{Y}(x) (\bar{H}_s^0 + \mathcal{Y}_s), \ \Lambda_n(x) = (\varsigma_0^2 + 1/\kappa)^{-1} \Lambda(x).$

Approximation of the state dynamics begins with an approximation of the input. The proof of Lemma 3.4 is postponed to Appendix C.

Lemma 3.4. The small-r approximation holds for the solution to IPD-Q:

$$\widetilde{\Phi}(u \mid s, r) = \Phi^{0}(u \mid s) \exp(\Lambda_{\mathsf{n}}(x)r) + O(r^{2}).$$

The following linear systems approximation follows easily from Lemma 3.4. We defer to [13] for details of the proof.

Proposition 3.5. Suppose that the input  $\phi_k(u \mid s) = \phi^*(u \mid s, r_k)$  is applied to the nonlinear system (1.22). The closed loop dynamics then admit the approximation

$$\widetilde{X}_{k+1} = A\widetilde{X}_k + BG^{\star}r_k + \varepsilon_k + O(r_k^2)$$

in which A and B are defined in Proposition 3.3,  $G^*$  is the column vector with entries  $G_i^{\star} = \Phi^0(u^i \mid s^i) \Lambda_{\mathsf{n}}(x^i)$ , and  $\varepsilon_k$  is quadratic in the deviation  $(\widetilde{\nu}_k, \Phi_k)$ :

$$\varepsilon_k^i = \sum_j \widetilde{\nu}_k(x^j) T(x^j, s^i) \widetilde{\Phi}_k(u^i \mid s^i) \,, \qquad x^i = (s^i, u^i) \in \mathsf{X}.$$

4. Applications to demand dispatch. An application of the control framework described in the previous sections is demand dispatch, an evolving science for automatically controlling flexible loads to help maintain supply-demand balance in the power grid. The goal of demand dispatch is to modify the behavior of flexible loads such that the aggregate power consumption tracks a reference signal that is broadcast by a BA.

Keep in mind that in the numerical examples here we focus entirely on the meanfield model. We know from prior work that evolution of the empirical distributions does closely track this idealization: for reasonably large  $\mathcal{N}$ , following the notation (1.1b), the approximation  $\nu_k^{\mathcal{N}} \approx \nu_k$  holds and the covariance of the error grows slowly with k (error is reduced with feedback [19, 20]). Although the control architecture in this prior work is very different, it should not surprise the reader that the law of large numbers and associated central limit theorem hold in the setting of this paper.

Also, the numerical results here focus entirely on the solutions surveyed in section 2. As explained in section 3, the IPD-Q solution for real-time feedback reduces to something similar to what has been extensively explored in prior work [19, 20].

Although these techniques can be applied to any flexible load, the experiments in this section demonstrate distributed control of a population of residential water heaters or refrigerators. An MDP model is constructed in which the state is the standard used to capture hysteresis control,  $S_k = (\theta_k, U_{k-1})$ , in which  $\theta_k \in \mathbb{R}$  is the temperature, and  $U_k \in \{0,1\}$  denotes power mode for each k. Remember the physical system operates in continuous time, and k represents the kth sampling time. This means that  $U_{k-1}$  represents the power mode during the sampling interval ending at the kth sampling time.

**4.1. Designing the nominal model.** Construction of the nominal model with transition matrices  $\{P_k^0\}$  of the form (1.10b) requires specification of dynamics of nature and the nominal policy. In the case of water heaters, the sequence of transition matrices  $\{T_k\}$  for nature were based on input-output data obtained from Oak Ridge National Laboratories [22]. For refrigerators, T was taken independent of k, constructed based on simulations of the standard linear TCL model,

(4.1) 
$$\theta_{k+1} = \theta_k + \alpha [\theta^a - \theta_k] - \beta U_k + D_{k+1},$$

in which  $\alpha, \beta > 0$ ,  $\theta^a$  denotes the (time-invariant) ambient air temperature, and the disturbance process D captures modeling error and usage.

In all cases the nominal policy was chosen time-homogeneous:  $\phi_k^0 \equiv \phi^0$  is a fixed randomized policy, designed to approximate deterministic hysteresis control. We describe the construction for water heaters, following [37, 22]. We defer to [13] for details on the construction of the nominal policy.

**4.2. Tracking.** In practical applications the aggregate power is of interest, which is approximated by  $\varrho \mathcal{N} y_k$  at time k, where  $\varrho$  is the rated power of a single load. Hence the total population size  $\mathcal{N}$  must be taken into account in any tracking problem. In plots that follow, we choose to focus on the "normalized" response, defined as follows:

$$y_k^{ ext{ref}} = r_k/\varrho \,, \qquad \qquad \hat{y}_k^{ ext{ref}} = \hat{r}_k/\varrho \,, \qquad \qquad y_k = \langle \nu_k, \mathcal{Y} \rangle/\varrho .$$

In this context,  $y_k$  can be interpreted as the probability of a load being on.

The two sets of plots in Figure 1 are distinguished by the reference signal. In each case the reference signal is a square wave. In (a) the signal is feasible, and in (b) it violates the energy limits of the collection of water heaters [27]. In Figure 1(a) it is seen that tracking is nearly perfect for sufficiently large  $\kappa$ . Tracking of the larger reference signal would require temperature deviations to exceed the deadband of the water heater. Instead, we observe in Figure 1(b) a graceful truncation of the reference signal.

The next set of experiments was designed to assess sensitivity of KLQ optimal control to modeling error. Specifically, what are the consequences of ignoring the randomness of nature?

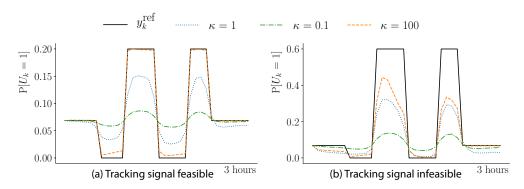


Fig. 1. Tracking error: (a) reference signal is feasible; (b) reference signal is infeasible.

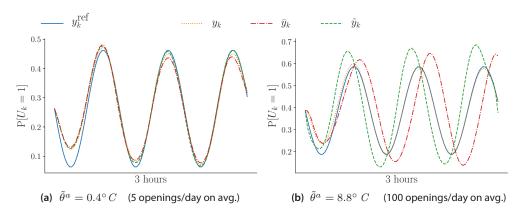


Fig. 2. Sensitivity experiments support the use of MPC with a deterministic approximation for randomness from nature.

A particular choice of statistics for (4.1) was chosen in order to mimic the effect of a refrigerator door opening at random times throughout the day: D is i.i.d., with

$$D_{k+1} = \begin{cases} \bar{d} & \text{with probability} \quad \varepsilon, \\ 0 & \text{with probability} \quad 1 - \varepsilon, \end{cases}$$

where  $\varepsilon$  determines the average amount of door openings per day, and  $\bar{d}$  was chosen so that the temperature inside the refrigerator increases when the door is open even when the power mode is on. A deterministic approximation of (4.1) was constructed for comparison, in which  $D_{k+1}$  is replaced by its mean:

(4.2) 
$$\theta_{k+1} = \theta_k + \alpha [\bar{\theta}^a - \theta_k] - \beta U_k$$

with  $\bar{\theta}^a = \theta^a + \tilde{\theta}^a$  with  $\tilde{\theta}^a = \bar{d}\varepsilon/\alpha$ .

Optimal policies were calculated for each of three models: the stochastic model (4.1), its deterministic approximation (4.2), and the cruder deterministic approximation obtained on setting  $D_{k+1} \equiv 0$  in (4.1) (equivalently, (4.2) with  $\bar{\theta}^a = \theta^a$ ). Each policy was then tested on the stochastic model (4.1).

Figure 2 displays the results from these experiments, where in each plot

- $y_k^{\text{ref}}$  is the reference signal,
- $y_k$  is the policy that is optimal for the stochastic model,

- $\tilde{y}_k$  is the policy that is optimal for (4.2),
- $\bar{y}_k$  is the policy that is optimal for (4.2) using  $\bar{\theta}^a = \theta^a$ .

The accurate tracking  $y_k \approx y_k^{\text{ref}}$  is expected because this reference signal is feasible, and  $\kappa > 0$  was chosen to be large.

It is seen in Figure 2(a) that all four trajectories are nearly identical for the smaller disturbance. The deviation is far greater in (b), for which the disturbance is greater. However,  $y_k$  and  $\bar{y}_k$  are nearly identical for about the first 30 minutes. This suggests that a deterministic approximation, combined with MPC, may be used in place of the stochastic model.

- **4.3. Information architectures.** The choice of information architecture is an interesting topic for future research. Here are three possibilities:
  - (i) Smart BA: The BA uses the reference signal  $\{r_k\}$  and its estimate of  $\nu_0^0$  to compute  $\lambda^*$  and broadcast it to the loads.
  - (ii) Smart load: The BA broadcasts  $\{r_k\}$  to the loads. Each load computes  $\lambda^*$  based on its internal model and  $\nu_0^0 = \delta_{x_0}$ , with  $x_0 \in X$  its current state.
  - (iii) Genius load: The BA broadcasts  $\{r_k\}$  to the loads. Each load computes  $\lambda^*$  based on its internal model and its estimate of  $\nu_0^0$ .

Each approach has its strengths and weaknesses. Approaches (i) and (iii) require knowledge of the initial marginal pmf of the population,  $\nu_0^0$ . If a perfect estimate is assumed, then the total cost in cases (i) and (iii) is equal to  $J^*(\nu_0^0)$ . But, how can a load estimate the marginal pmf of the population? Numerical results from [13] suggest coupling of the marginals from distinct initial conditions. If enough time has passed since the latest MPC iteration, the pmfs  $\{\nu_k\}$  computed locally can be used to approximate the marginal pmf of the population (perhaps smoothed using the techniques of [19, 20]).

In contrast, the total cost for case (ii) is the sum,  $\sum_{i=1}^{d} \nu_0^0(x^i) J^*(\delta_{x^i})$ , since each load optimizes according to its own initial state,  $x^i$ . Even when the aggregate can easily track  $\{r_k\}$ , the cost  $J^*(\delta_{x^i})$  may be very large for individuals that are at odds with the reference signal. For example, an increase in power consumption could be requested while a water heater is near its upper temperature limit and must turn off. So, it is possible that approach (ii) will impose greater stress on the loads as compared to the other two options, or will lead to reduced capacity.

- **5.** Conclusions. The paper provides a complete theory for KLQ and infinite-horizon counterparts, without the restriction to deterministic dynamics imposed in [15, 23]. Plans for future research include the following:
  - (i) Monte Carlo approaches for both KLQ and IPD-Q. The approximation (1.27) invites actor critic methods for approximating the best coefficients  $\{K_{i,j}^{\star}, G_{i}^{\star}\}$  based on training data with nonconstant reference signal, rather than approximation.
  - (ii) Evaluate robustness and sensitivity to other types of modeling error.
  - (iii) Investigate alternative transform techniques.
  - (iv) Consider other cost functions, such as the Wasserstein distance. Preliminary results are summarized in [24].
  - (v) Investigate the relationship between optimality and coupling of the pmfs, and the implications to control design.
  - (vi) Careful design of a terminal cost function may result in better performance for smaller time horizons [18].
  - (vii) How is the relative value function  $\mathcal{H}$  appearing in (3.1) related to  $h^{\lambda}$  appearing in (3.5) (with  $\lambda = \lambda_r$ )?

**Appendix A. Convexity.** The proofs of Theorem 2.1 and Proposition 2.2 make use of the following four lemmas. The first is based on a well-known result regarding relative entropy. For any function  $h: \mathsf{X}^{K+1} \to \mathbb{R}$  denote  $\Lambda^0(h) := \sup_p \big\{ \langle p, h \rangle - D(p \| p^0) \big\}$ .

LEMMA A.1 (convex dual of relative entropy). For each  $p^0 \in \mathcal{S}(\mathsf{X}^{K+1})$  and function  $h \colon \mathsf{X}^{K+1} \to \mathbb{R}$ , the (possibly infinite) value of  $\Lambda^0(h)$  coincides with the log moment generating function:  $\Lambda^0(h) = \log \langle p^0, e^h \rangle$ . Moreover, provided  $\Lambda^0(h) < \infty$ , the supremum defining this quantity is uniquely attained with  $p^* = p^0 \exp(h - \Lambda^0(h))$ . That is, the log-likelihood  $L^* = \log(dp^*/dp^0)$  is given by  $L^*(\vec{x}) = h(\vec{x}) - \Lambda^0(h)$ .

Lemma A.2. The dual function can be expressed

(A.1) 
$$\varphi^{\star}(\lambda, g) = \lambda^{T} \hat{r} - \frac{1}{2\kappa} \|\lambda\|^{2} - \langle \nu_{0}^{0}, G_{1}^{\lambda} \rangle + \sum_{k=1}^{K} \min_{s} \left[ g_{k}(s) - \mathcal{T}_{k}^{\lambda}(g_{k+1}; s) \right].$$

*Proof.* First, make the substitution  $\nu_k(s,u) = \hat{\nu}_k(s) \phi_k(u \mid s)$ , so that the Lagrangian (2.3) can be written

$$\mathcal{L}(\nu, \gamma, \lambda, g) = \sum_{n=1}^{N} \left( \frac{\kappa}{2} \gamma_n^2 + \lambda_n \gamma_n + \lambda_n \hat{r}_n \right) - \sum_{s,u} \nu_0^0(s, u) \sum_{s'} T_0(x, s') g_1(s')$$

$$+ \sum_{k=1}^{K} \sum_{s} \hat{\nu}_k(s) \sum_{u} \phi_k(u \mid s) \left( L_k(s, u) - \sum_{s'} T_k(x, s') g_{k+1}(s') - \check{\lambda}_k \mathcal{Y}(s, u) \right)$$

$$+ \sum_{k=1}^{K} \sum_{s} \hat{\nu}_k(s) g_k(s)$$

with  $g_{K+1} \equiv 0$ , and  $L_k(s,u) = \log \frac{\phi_k(u|s)}{\phi_k^0(u|s)}$ . This amounts to a Lagrangian decomposition since the minimization of the Lagrangian is equivalent to solving K separate convex programs to obtain each of the minimizers  $\{\nu_k^{\lambda,g}: \nu_k^{\lambda,g}(s,u) = \hat{\nu}_k^{\lambda,g}(s)\phi_k^{\lambda,g}(u \mid s), (s,u) \in \mathsf{X}, 1 \leq k \leq K\}$ . That is,  $\arg \min_{\Phi} \mathcal{L} =$ 

(A.3) 
$$\left\{ \underset{\phi_k: 1 \leq k \leq K}{\operatorname{arg \, min}} \sum_{u} \phi_k(u \mid s) \left[ L_k(s, u) - \sum_{s'} T_k(x, s') g_{k+1}(s') - \check{\lambda}_k \mathcal{Y}(s, u) \right] \right\}.$$

Lemma A.1 implies that the minimizer is obtained with  $\Lambda_k(s) = \mathcal{T}_k^{\lambda}(g_{k+1};s)$  and

$$(A.4) \qquad \Phi_k^{\lambda,g}(u \mid s) = \Phi_k^0(u \mid s) \exp\left(\sum_{s'} T_k(x, s') g_{k+1}(s') + \check{\lambda}_k \mathcal{Y}(s, u) - \Lambda_k(s)\right).$$

Lemma A.1 also gives the value

$$\min_{\Phi_k} \sum_{u} \Phi_k(u \mid s) \left[ L_k(s, u) - \sum_{s'} T_k(x, s') g_{k+1}(s') - \check{\lambda}_k \mathcal{Y}(s, u) \right] = -\mathcal{T}_k^{\lambda}(g_{k+1}; s)$$

resulting in

(A.5) 
$$\min_{\nu} \mathcal{L}(\nu, \gamma, \lambda, g) = \sum_{n=1}^{N} \left( \frac{\kappa}{2} \gamma_n^2 + \lambda_n \gamma_n + \lambda_n \hat{r}_n \right) - \sum_{s,u} \nu_0^0(s, u) \sum_{s'} T_0(x, s') g_1(s') + \sum_{k=1}^{K} \min_{\hat{\nu}_k} \langle \hat{\nu}_k, g_k - \mathcal{T}_k^{\lambda}(g_{k+1}) \rangle.$$

Next, observe that the minimizer  $\hat{\nu}_k^{\lambda,g}$  is obtained when the support of each  $\hat{\nu}_k$  satisfies

$$\operatorname{supp}(\hat{\nu}_k(s)) \subseteq \underset{s}{\operatorname{arg\,min}} \left[ g_k(s) - \mathcal{T}_k^{\lambda}(g_{k+1};s) \right]$$
so that 
$$\min_{s} \left[ g_k(s) - \mathcal{T}_k^{\lambda}(g_{k+1};s) \right] = \langle \hat{\nu}_k^{\lambda,g}, g_k - \mathcal{T}_k^{\lambda}(g_{k+1}) \rangle.$$

We also have  $\gamma_n^{\lambda} = -\frac{1}{\kappa}\lambda_n$  Substituting the minimizers  $\{\nu_k^{\lambda,g}, \gamma_n^{\lambda}\}$  into (A.5), and applying (2.7), results in (A.1).

## Appendix B. Duality.

Lemma B.1. The maximum of the dual function over g is

(B.1) 
$$\varphi^{\star}(\lambda) := \max_{g} \varphi^{\star}(\lambda, g) = \lambda^{T} \hat{r} - \frac{1}{2\kappa} \|\lambda\|^{2} - \langle \nu_{0}^{0}, G_{1}^{\lambda} \rangle$$

with  $G_1^{\lambda}(x) = \sum_{s'} T_0(x, s') g_1^{\lambda}(s')$ . A maximizer  $g^{\lambda}$  is given by the recursive formula:

(B.2) 
$$g_k^{\lambda} = \mathcal{T}_k^{\lambda}(g_{k+1}^{\lambda}), \ 1 \le k \le K, \ where \ g_{K+1}^{\lambda} \equiv 0.$$

*Proof.* Adding a constant to any of the  $(g_1, g_2, \ldots, g_K)$  does not change the value of  $\mathcal{L}$  (2.3) or  $\varphi^*$  (2.8), so without loss of generality we assume

(B.3) 
$$\min_{s} \left[ g_k(s) - \mathcal{T}_k^{\lambda}(g_{k+1}; s) \right] = 0 \text{ for each } k,$$

and consequently

(B.4) 
$$g_k \ge \mathcal{T}_k^{\lambda}(g_{k+1})$$
 for each  $k$ .

Thus, in view of (A.1),

(B.5) 
$$\varphi^{\star}(\lambda) = \lambda^{T} \hat{r} - \frac{1}{2\kappa} \|\lambda\|^{2} - \min_{g_{1}} \sum_{s,u} \nu_{0}^{0}(s,u) \sum_{s'} T_{0}(x,s') g_{1}(s'),$$

where the minimum is subject to the constraint (B.4). Next, observe that  $\mathcal{T}_k^{\lambda}$  is a monotone operator, so that for each  $k \leq K$ ,

$$g_k \ge \mathcal{T}_k^{\lambda} \circ \mathcal{T}_{k+1}^{\lambda} \circ \cdots \circ \mathcal{T}_K^{\lambda}(g_{K+1}) \doteq g_k^{\lambda}, \text{ where } g_{K+1} \equiv 0.$$

Based on the expression (B.5), we now show that the maximum  $\arg \max_g \phi^*(\lambda, g)$  is obtained by choosing each  $g_k$  to reach this lower bound, giving (B.2). Indeed,  $g_1^{\lambda}$  achieves the minimum in (B.5), since  $g_1^{\lambda} \leq g_1$  for any  $g_1$  for which (B.4) holds. This result along with (B.3) yields (B.1).

For an inductive proof of the following see [13].

LEMMA B.2. The maximizers  $\{g_k^{\lambda}\}$  have at most linear growth in  $\|\lambda\|$ :

(B.6) 
$$|g_k^{\lambda}(s)| \le C_k ||\lambda||, \quad 1 \le k \le K,$$

where  $C_k = \|\mathcal{Y}\|_{\infty} \sum_{i=k}^{K} \|w(i)\|$  and w(i) is the vector  $\{w_1(i), w_2(i), \dots, w_N(i)\}$ .

Proof of Theorem 2.1. We prove the existence of a maximizer  $\lambda^*$  by showing that  $\phi^*(\lambda)$  is an anticoercive function, i.e.,  $\phi^*(\lambda) \to -\infty$  as  $\|\lambda\| \to \infty$ . By Lemma B.2, there exists  $C_1 < \infty$  such that

$$\varphi^{*}(\lambda) = \lambda^{T} \hat{r} - \frac{1}{2\kappa} \|\lambda\|^{2} - \sum_{s,u} \nu_{0}^{0}(s,u) \sum_{s'} T_{0}(x,s') g_{1}^{\lambda}(s')$$

$$\leq \|\lambda\| \|\hat{r}\| - \frac{1}{2\kappa} \|\lambda\|^{2} + \max_{s'} |g_{1}^{\lambda}(s')| \leq \|\lambda\| \|\hat{r}\| - \frac{1}{2\kappa} \|\lambda\|^{2} + C_{1} \|\lambda\|.$$

Since  $\phi^*(\lambda)$  is upper-bounded by an anticoercive function,  $\phi^*(\lambda)$  itself is an anticoercive function. Thus a maximizer  $\lambda^*$  exists, and  $(\lambda^*, g^*) = (\lambda^*, g^{\lambda^*})$  by (B.2).

The primal is a convex program, as established in Proposition 1.1. To show that there is no duality gap it is sufficient that Slater's condition holds [6, section 5.3.2]. This condition holds: the relative interior of the constraint set for the primal is nonempty since it contains  $\{\nu_k^0\}$ . Optimality of (2.6) is established by substituting  $g_{k+1}^*$  into (A.4) and by making the substitution  $g_k^* = \mathcal{T}_k^{\lambda}(g_{k+1}^*)$  implied by (B.2).

Proof of Proposition 2.2. This proof has three parts:

- (i) Equation (2.5) is proven by Lemma B.1.
- (ii) Equation (2.8) is proven by Lemma B.1.
- (iii) The representation of the derivative in part (iii) is standard (e.g., section 5.6 of [6], or [13]).

# Appendix C. IPD-Q.

Proof of Lemma 3.4. An application of the implicit function theorem tells us that  $\{r^{\lambda}, \phi^{\lambda} : \lambda \in \mathbb{R}\}$  are smooth as functions of  $\lambda$ , whose derivatives may be expressed

$$\frac{d}{d\lambda}r^{\lambda} = \frac{d}{d\lambda}\langle \pi^{\lambda}, \mathcal{Y} \rangle + 1/\kappa = \varsigma_{\lambda}^{2} + 1/\kappa,$$
$$\frac{d}{d\lambda}\log(\phi^{\lambda}(u \mid s)) = \bar{H}^{\lambda}(x) + \mathcal{Y}(x) - \frac{d}{d\lambda}\Gamma^{\lambda}(s).$$

The first identities follow from (3.4) and then (3.8). The formula for the derivative of  $\log(\phi^{\lambda})$  is immediate from (3.6).

The proof of (3.11) requires approximations for  $r^{\lambda}$  and  $\phi^{\lambda}$  in a neighborhood of zero. The first approximation is given by  $r^{\lambda} = (\varsigma_0^2 + 1/\kappa)\lambda + O(\lambda^2)$ . The definition (3.7) implies that

$$\frac{d}{d\lambda} \Gamma^{\lambda}(s)\big|_{\lambda=0} = \bar{H}_s^0 + \mathcal{Y}_s,$$
 which gives 
$$\log \left( \Phi^{\lambda}(u \mid s) \right) = \log \left( \Phi^0(u \mid s) \right) + \Lambda(x)\lambda + O(\lambda^2).$$

An inversion is applied to express  $\lambda$  as a function of r, giving

$$\phi(u \mid s, r) = \phi^0(u \mid s) \exp(\Lambda(x)\lambda_r) + O(r^2)$$

with  $\lambda_r = (\varsigma_0^2 + 1/\kappa)^{-1} r + O(r^2)$ . Hence (3.11) follows from a first-order Taylor series approximation of the exponential.

### REFERENCES

- [1] M. Almassalkhi, J. Frolik, and P. Hines, *Packetized energy management: Asynchronous and anonymous coordination of thermostatically controlled loads*, in Proceedings of the American Control Conference, 2017, pp. 1431–1437.
- [2] E. Altman, Constrained Markov Decision Processes, Chapman & Hall/CRC, Boca Raton, FL, 1999.
- [3] B. D. O. Anderson and J. B. Moore, Optimal Control: Linear Quadratic Methods, Prentice-Hall, Englewood Cliffs, NJ, 1990.
- [4] E. BENENATI, M. COLOMBINO, AND E. DALL'ANESE, A tractable formulation for multi-period linearized optimal power flow in presence of thermostatically controlled loads, in Proceedings of the IEEE Conference on Decision and Control, 2019, pp. 4189–4194.
- [5] V. S. BORKAR, Convex analytic methods in Markov decision processes, in Handbook of Markov Decision Processes, Internat. Ser. Oper. Res. Management Sci. 40, Kluwer, Boston, 2002, pp. 347–375.

- [6] S. BOYD AND L. VANDENBERGHE, Convex Optimization, Cambridge University Press, New York, 2004.
- [7] A. BROOKS, E. LU, D. REICHER, C. SPIRAKIS, AND B. WEIHL, Demand dispatch, IEEE Power Energy Mag., 8 (2010), pp. 20–29.
- [8] A. Bušić and S. Meyn, Distributed randomized control for demand dispatch, in Proceedings of the Conference on Decision and Control, 2016, pp. 6964–6971.
- [9] A. Bušić and S. Meyn, Ordinary differential equation methods for Markov decision processes and application to Kullback-Leibler control cost, SIAM J. Control Optim., 56 (2018), pp. 343–366.
- [10] P. E. CAINES, Mean field games, in Encyclopedia of Systems and Control, J. Baillieul and T. Samad, eds., Springer, London, 2021, pp. 1197–1202.
- [11] A. Bušić, S. Meyn, and N. Cammardella, Learning optimal policies in mean-field models with Kullback-Leibler regularization, in Proceedings of the IEEE Conference on Decision and Control, 2023.
- [12] N. CAMMARDELLA, Creating Virtual Energy Storage through Optimal Allocation and Control of Flexible Power Consumption, Ph.D. thesis, University of Florida, Gainesville, 2021.
- [13] N. Cammardella, A. Bušić, and S. Meyn, Kullback-Leibler-Quadratic Optimal Control, https://arxiv.org/abs/2004.01798, 2020.
- [14] N. CAMMARDELLA, A. Bušić, And S. Meyn, Simultaneous allocation and control of distributed energy resources via Kullback-Leibler-quadratic optimal control, in Proceedings of the American Control Conference, 2020, pp. 514–520.
- [15] N. CAMMARDELLA, A. BUŠIĆ, Y. JI, AND S. MEYN, Kullback-Leibler-quadratic optimal control of flexible power demand, in Proceedings of the Conference on Decision and Control, 2019, pp. 4195–4201.
- [16] R. CARMONA AND F. DELARUE, Probabilistic Theory of Mean Field Games with Applications I: Mean Field FBSDEs, Control, and Games, Probability Theory and Stochastic Modelling, Springer, New York, 2018.
- [17] R. CARMONA AND F. DELARUE, Probabilistic Theory of Mean Field Games with Applications II: Mean Field Games with Common Noise and Master Equations, Probability Theory and Stochastic Modelling, Springer, New York, 2018.
- [18] R.-R. CHEN AND S. P. MEYN, Value iteration and optimization of multiclass queueing networks, Queueing Syst., 32 (1999), pp. 65–97.
- [19] Y. CHEN, Markovian Demand Dispatch Design for Virtual Energy Storage to Support Renewable Energy Integration, Ph.D. thesis, University of Florida, Gainesville, 2016.
- [20] Y. CHEN, A. BUŠIĆ, AND S. MEYN, State estimation for the individual and the population in mean-field control with application to demand dispatch, IEEE Trans. Automat. Control, 62 (2017), pp. 1138–1149.
- [21] Y. CHEN, T. T. GEORGIOU, AND M. PAVON, Optimal transport in systems and control, Annu. Rev. Control Robot. Auton. Syst., 4 (2020), pp. 89–113.
- [22] Y. CHEN, M. U. HASHMI, J. MATHIAS, A. BUŠIĆ, AND S. MEYN, Distributed control design for balancing the grid using flexible loads, in Energy Markets and Responsive Grids: Modeling, Control, and Optimization, S. Meyn, T. Samad, I. Hiskens, and J. Stoustrup, eds., Springer, New York, 2018, pp. 383–411.
- [23] M. CHERTKOV AND V. Y. CHERNYAK, Ensemble control of cycling energy loads: Markov decision approach, in Energy Markets and Responsive Grids: Modeling, Control, and Optimization, IMA Vol. Control Energy Markets Grids 162, Springer, New York, 2018.
- [24] T. L. CORRE AND S. M. ANA BUŠIĆ, Feature Projection for Optimal Transport, https://arxiv.org/abs/2208.01958, 2023.
- [25] A. Dembo and O. Zeitouni, Large Deviations Techniques and Applications, 2nd ed., Springer, New York, 1998.
- [26] O. Guéant, J.-M. Lasry, and P.-L. Lions, Mean Field Games and Applications, Springer, Berlin, 2011, pp. 205–266.
- [27] H. HAO, B. M. SANANDAJI, K. POOLLA, AND T. L. VINCENT, Aggregate flexibility of thermostatically controlled loads, IEEE Trans. Power Syst., 30 (2015), pp. 189–198.
- [28] M. HUANG, P. E. CAINES, AND R. P. MALHAME, Large-population cost-coupled LQG problems with nonuniform agents: Individual-mass behavior and decentralized ε-Nash equilibria, IEEE Trans. Automat. Control, 52 (2007), pp. 1560–1571.
- [29] M. Huang, R. P. Malhame, and P. E. Caines, Large population stochastic dynamic games: Closed-loop Mckean-Vlasov systems and the Nash certainty equivalence principle, Commun. Inf. Syst., 6 (2006), pp. 221–251.
- [30] I. KONTOYIANNIS AND S. P. MEYN, Large deviations asymptotics and the spectral theory of multiplicatively regular Markov processes, Electron. J. Probab., 10 (2005), pp. 61–123.

- [31] J. M. LASRY AND P. L. LIONS, Mean field games, Jpn. J. Math., 2 (2007), pp. 229–260.
- [32] J.-S. Li, Ensemble control of finite-dimensional time-varying linear systems, IEEE Trans. Automat. Control, 56 (2010), pp. 345–357.
- [33] D. G. LUENBERGER, Optimization by Vector Space Methods, John Wiley & Sons, New York, 1969, reprinted 1997.
- [34] A. S. Manne, Linear programming and sequential decisions, Management Sci., 6 (1960), pp. 259–267.
- [35] J. MATHIEU, S. KOCH, AND D. CALLAWAY, State estimation and control of electric loads to manage real-time energy imbalance, IEEE Trans. Power Syst., 28 (2013), pp. 430–440.
- [36] S. Meyn, Control Systems and Reinforcement Learning, Cambridge University Press, Cambridge, UK, 2022.
- [37] S. MEYN, P. BAROOAH, A. BUŠIĆ, Y. CHEN, AND J. EHREN, Ancillary service to the grid using intelligent deferrable loads, IEEE Trans. Automat. Control, 60 (2015), pp. 2847–2862.
- [38] N. PARIKH AND S. BOYD, Proximal Algorithms, Foundations and Trends in Optimization, Now Publishers, 2013.
- [39] G. Peyré and M. Cuturi, Computational Optimal Transport, https://arxiv.org/abs/1803.00567, 2020.
- [40] J. C. Principe, Information Theoretic Learning: Renyi's Entropy and Kernel Perspectives, Springer, New York, 2010.
- [41] M. L. PUTERMAN, Markov Decision Processes: Discrete Stochastic Dynamic Programming, John Wiley & Sons, New York, 2014.
- [42] A. Taghvaei and P. G. Mehta, A survey of feedback particle filter and related controlled interacting particle systems (CIPS), Annu. Rev. Control, 55 (2023), pp. 356–378, https://doi.org/10.1016/j.arcontrol.2023.03.006.
- [43] S. H. TINDEMANS, V. TROVATO, AND G. STRBAC, Decentralized control of thermostatic loads for flexible demand response, IEEE Trans. Control Syst. Technol., 23 (2015), pp. 1685–1700.
- [44] E. TODOROV, Linearly-solvable Markov decision problems, in Proceedings of Advances in Neural Information Processing Systems, B. Schölkopf, J. Platt, and T. Hoffman, eds., Cambridge, MA, 2007, pp. 1369–1376.