Preliminary Results on Control and Random Dynamical Systems in Reproducing Kernel Hilbert Spaces

> Boumediene Hamzi, Department of Mathematics, Imperial College, London, UK. email: b.hamzi@imperial.ac.uk

joint work with Jake Bouvrie (Duke University)

Goal: Combining tools from the theories of Dynamical Systems and Learning in view of analyzing, predicting and controling nonlinear systems on the basis of data rather than models. Approach: View Reproducing Kernel Hilbert Spaces as "Linearizing Spaces", i.e. Nonlinear Systems will be mapped into an RKHS where Linear Systems Theory will be applied.

- Review of Some Concepts for Linear Control Systems
- Elements of Learning Theory
- On Nonlinear Control Systems in RKHS
- Application: Model Reduction of Nonlinear Control Systems in RKHS
- Review of Some Concepts for Linear SDEs
- On Nonlinear SDEs in RKHSes
- Application: Estimation of the Stationary Solution of the Fokker-Planck Equation of nonlinear SDEs
- Application: Parameter Estimation of SDEs in RKHSes

• Consider a linear control system

$$\begin{array}{rcl} \dot{x} &=& Ax + Bu \\ y &=& Cx \end{array},$$

where $x \in \mathbb{R}^n$, $u \in \mathbb{R}^q$, $y \in \mathbb{R}^p$, (A, B) is controllable, (A, C) is observable and A is Hurwitz.

- We define the controllability and the observability Gramians as, respectively, $W_c = \int_0^\infty e^{At} B B^{\mathsf{T}} e^{A^{\mathsf{T}}t} dt$, $W_o = \int_0^\infty e^{A^{\mathsf{T}}t} C^{\mathsf{T}} C e^{At} dt$.
- These two matrices can be viewed as a measure of the controllability and the observability of the system.

• Consider the past energy, $L_c(x_0)$, defined as the minimal energy required to reach x_0 from 0 in infinite time

$$L_c(x_0) = \inf_{\substack{u \in L_2(-\infty,0), \\ x(-\infty) = 0, x(0) = x_0}} \frac{1}{2} \int_{-\infty}^0 \|u(t)\|^2 \, dt.$$

• Consider the future energy, $L_o(x_0)$, defined as the output energy generated by releasing the system from its initial state $x(t_0) = x_0$, and zero input u(t) = 0 for $t \ge 0$, i.e.

$$L_o(x_0) = \frac{1}{2} \int_0^\infty \|y(t)\|^2 \, dt,$$

for $x(t_0) = x_0$ and $u(t) = 0, t \ge 0$.

• In the linear case, it can be shown that

$$L_c(x_0) = \frac{1}{2} x_0^{\mathsf{T}} W_c^{-1} x_0, \quad L_o(x_0) = \frac{1}{2} x_0^{\mathsf{T}} W_o x_0.$$

 \bullet Moreover, W_c and W_o satisfy the following Lyapunov equations

 $AW_c + W_c A^{\mathsf{T}} = -BB^{\mathsf{T}}, \quad A^{\mathsf{T}}W_o + W_o A = -C^{\mathsf{T}}C.$

- These energies are directly related to the controllability and observability operators.
- Given a matrix pair (A, B), the controllability operator Ψ_c is defined as $\Psi_c: L_2(-\infty, 0) \to \mathbb{C}^n; u \mapsto \int_{-\infty}^0 e^{-A\tau} Bu(\tau) d\tau.$
- The significance of this operator is made evident via the following optimal control problem: Given the linear system $\dot{x}(t) = Ax(t) + Bu(t)$ defined for $t \in (-\infty, 0)$ with $x(-\infty) = 0$, and for $x(0) \in \mathbb{C}^n$ with unit norm, what is the minimum energy input u which drives the state x(t) to $x(0) = x_0$ at time zero? That is, what is the $u \in L_2(-\infty, 0]$ solving $\Psi_c u = x_0$ with smallest norm $||u||_2$?

- •If (A, B) is controllable, then $\Psi_c \Psi_c^* =: W_c$ is nonsingular, and the answer to the preceding question is $u_{opt} := \Psi_c^* W_c^{-1} x_0$. The input energy is given by $\|u_{opt}\|_2^2 = x_0^* W_c^{-1} x_0$.
- The reachable set through u_{opt} , i.e. the final states $x_0 = \Psi_c u$ that can be reached given an input $u \in L_2(-\infty, 0]$ of unit norm, $\{\Psi_c u : u \in L_2(-\infty, 0] \text{ and } \|u\|_2 \leq 1\}$ may be defined as

$$\mathcal{R} := \{ W_c^{\frac{1}{2}} x_c : x_c \in \mathbb{C}^n \text{ and } \|x_c\| \le 1 \}.$$

Review of Some Concepts for Linear Control Systems

• For the autonomous system $\dot{x} = Ax, x(0) = x_0 \in \mathbb{C}^n; y = Cx$ where A is Hurwitz, the observability operator is defined as

$$\Psi_o: \mathbb{C}^n \to L_2(0,\infty); x_0 \mapsto \begin{cases} Ce^{At}x_0, & \text{for } t \ge 0\\ 0, & \text{otherwise} \end{cases}$$

• The corresponding observability ellipsoid is given by

$$\mathcal{E} := \{ W_o^{\frac{1}{2}} x_0 : x_0 \in \mathbb{C}^n \text{ and } \|x_0\| = 1 \}.$$

• The energy of the output signal $y = \Psi_o x_0$, for $x_0 \in \mathbb{C}^n$ can then be computed as

$$\|y\|_2^2 = \langle \Psi_o x_0, \Psi_o x_0 \rangle = \langle x_0, \Psi_o^* \Psi_o x_0 \rangle = \langle x_0, W_0 x_0 \rangle$$

where $\Psi_o^*: L_2[0,\infty) \to \mathbb{C}^n$ is the adjoint of Ψ_o .

 \bullet Consider the nonlinear system Σ

$$\begin{cases} \dot{x} = f(x) + \sum_{i=1}^{m} g_i(x)u_i, \\ y = h(x), \end{cases}$$

with $x \in \mathbb{R}^n$, $u \in \mathbb{R}^m$, $y \in \mathbb{R}^p$, f(0) = 0, $g_i(0) = 0$ for $1 \le i \le m$, and h(0) = 0.

Hypothesis H: The linearization of around the origin is controllable, observable and $F = \frac{\partial f}{\partial x}|_{x=0}$ is asymptotically stable.

Controllability and Observability Energies for Nonlinear System

• Theorem (Scherpen, 1993) If the origin is an asymptotically stable equilibrium of f(x) on a neighborhood W of the origin, then for all $x \in W$, $L_o(x)$ is the unique smooth solution of

$$\frac{\partial L_o}{\partial x}(x)f(x) + \frac{1}{2}h^{\mathsf{T}}(x)h(x) = 0, \quad L_o(0) = 0$$

under the assumption that this equation has a smooth solution on W. Furthermore for all $x \in W$, $L_c(x)$ is the unique smooth solution of

$$\frac{\partial L_c}{\partial x}(x)f(x) + \frac{1}{2}\frac{\partial L_c}{\partial x}(x)g(x)g^{\mathsf{T}}(x)\frac{\partial^{\mathsf{T}}L_c}{\partial x}(x) = 0, \quad L_c(0) = 0$$

under the assumption that this equation has a smooth solution L_c on Wand that the origin is an asymptotically stable equilibrium of $-(f(x) + g(x)g^{\mathsf{T}}(x)\frac{\partial \bar{L}_c}{\partial x}(x))$ on W.

Boumediene Hamzi (Imperial College)

On Control and RDS in RKHS

June 4th, 2012 11 / 55

Controllability and Observability Energies in Model Reduction of Linear Control Systems

• Gramians have several uses in Linear Control Theory. For example, for the purpose of model reduction.

• Balancing: find a representation where the system's observable and controllable subspaces are aligned so that reduction, if possible, consists of eliminating uncontrollable states which are also the least observable.

• More formally, we would like to find a new coordinate system such that

$$W_c = W_o = \Sigma = \operatorname{diag}\{\sigma_1, \cdots, \sigma_n\},\$$

where $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n \geq 0$. If (F, G) is controllable and (F, H) is observable, then there exists a transformation such that the state space expressed in the transformed coordinates (TFT^{-1}, TG, HT^{-1}) is balanced and $TW_cT^{\top} = T^{-\top}W_oT^{-1} = \Sigma$.

Typically one looks for a gap in the singular values {σ_i} for guidance as to where truncation should occur. If we see that there is a k such that σ_k ≫ σ_{k+1}, then the states most responsible for governing the input-output relationship of the system are (x₁, ..., x_k) while (x_{k+1},..., x_n) are assumed to make negligible contributions.
Although several methods exist for computing T, the general idea is to compute the Cholesky decomposition of W_o so that W_o = ZZ^T, and form the SVD UΣ²U^T of Z^TW_cZ. Then T is given by T = Σ^{1/2}U^TZ⁻¹.

• Theorem (Scherpen) Consider system Σ under Hypothesis H and the assumptions in the preceding theorem. Then, there exists a neighborhood W of the origin and coordinate transformation $x = \varphi(z)$ on W converting the energy functions into the form

$$L_c(\varphi(z)) = \frac{1}{2} z^{\mathsf{T}} z,$$

$$L_o(\varphi(z)) = \frac{1}{2} \sum_{i=1}^n z_i^2 \sigma_i(z_i)^2,$$

where $\sigma_1(x) \ge \sigma_2(x) \ge \cdots \ge \sigma_n(x)$. The functions $\sigma_i(\cdot)$ are called *Hankel* singular value functions.

- In the above framework for balancing of nonlinear systems, one needs to solve (or numerically evaluate) the PDEs and compute the coordinate change $x = \varphi(z)$.
- However there are no systematic methods or tools for solving these equations.
- Various approximate solutions based on Taylor series expansions have been proposed Krener (2007, 2008), Fujimoto and Tsubakino (2008).
- Newman and Krishnaprasad (2000) introduce a statistical approximation based on exciting the system with white Gaussian noise and then computing the balancing transformation using an algorithm from differential topology.
- An essentially linear empirical approach, similar to Moore's empirical approach, was proposed by Lall, Marsden and Glavaski (2002).

Computing the Controllability and Observability Energies: Linear Case

• Analytic Approach: The Gramians W_c and W_o satisfy the Lyapunov equations

$$FW_c + W_c F^{\mathsf{T}} = -GG^{\mathsf{T}},$$

$$F^{\mathsf{T}}W_o + W_o F = -H^{\mathsf{T}}H.$$

• Data-Based Approach: Moore showed that W_c and W_o can be obtained from the impulse responses of Σ_L . For instance,

$$W_c = \int_0^\infty X(t)X(t)^T dt, \quad W_o = \int_0^\infty Y^T(t)Y(t)dt$$

where X(t) is the response to $u^i(t) = \delta(t)e_i$ with x(0) = 0, and Y(t) is the output response to u(t) = 0 and $x(0) = e_i$. Given X(t) and Y(t), one can perform PCA to obtain W_c and W_o respectively. The observability and controllability Gramians may be estimated statistically from typical system trajectories:

$$\widehat{W}_c = \frac{T}{mN} \sum_{i=1}^N X(t_i) X(t_i)^{\mathsf{T}}, \quad \widehat{W}_o = \frac{T}{pN} \sum_{i=1}^N Y(t_i) Y(t_i)^{\mathsf{T}}.$$

where $t_i \in [0, T], i = 1, ..., N$, $X(t) = [x^1(t) \cdots x^m(t)]$, and $Y(t) = [y^1(t) \cdots y^n(t)]^\top$ if $\{x^j(t)\}_{j=1}^m, \{y^j(t)\}_{j=1}^n$ are measured (vector-valued) responses and outputs of the system.

Computing the Controllability and Observability Energies for Nonlinear Systems

Questions

• How to extend Moore's empirical approach to Nonlinear Control Systems ?

• Are there "Gramians" for Nonlinear Systems ? and in the affirmative, how to compute them from data ?

• Can we "view" a Nonlinear (Control) System as Linear by working in a different space or at least perform PCA for Nonlinear Systems ?

• Idea ! Use of kernel methods. A kernel based procedure may be interpreted as mapping the data, through "feature maps", from the original input space into a potentially higher dimensional Reproducing Kernel Hilbert Space where linear methods may then be used.

- Historical Context: Appeared in the 1930s as an answer to the question: when is it possible to embed a metric space into a Hilbert space ? (Schoenberg, 1937)
- Answer: If the metric satisfies certain conditions, it is possible to embed a metric space into a special type of Hilbert spaces called RKHSes.
- Properties of RKHSes have been further studied in the 1950s and later (Aronszajn, 1950; Schwartz, 1964 etc.)

- Definition: A Hilbert Space is an inner product space that is complete and separable with respect to the norm defined by the inner product.
- Definition: For a compact $\mathcal{X} \subseteq \mathbb{R}^d$, and a Hilbert space \mathcal{H} of functions $f : \mathcal{X} \to \mathbb{R}$, we say that \mathcal{H} is a RKHS if there exists $k : \mathcal{X} \to \mathbb{R}$ such that
 - i. k has the reproducing property, i.e. $\forall f \in \mathcal{H}, f(x) = \langle f(\cdot), k(\cdot, x) \rangle.$
 - ii. k spans \mathcal{H} , i.e. $\mathcal{H} = \overline{\text{span}\{k(x, \cdot) | x \in \mathcal{X}\}}$.

• Definition: A Reproducing Kernel Hilbert Space (RKHS) is a Hilbert space H with a reproducing kernel whose span is dense in H. Equivalently, an RKHS is a Hilbert space of functions with all evaluation functionals bounded and linear.

• Remark: L_2 is a Hilbert space but not an RKHS because the delta function which has the reproducing property $f(x) = \int \delta(x-u)f(u)du$, does not satisfy the square integrable condition $\int \delta(u)^2 du \not< \infty$

The important properties of reproducing kernels are

- K(x,y) is unique.
- $\forall x, y \in \mathcal{X}$, K(x, y) = K(y, x) (symmetry).
- $\sum_{i,j=1}^{m} \alpha_i \alpha_j K(x_i, x_j) \ge 0$ for $\alpha_i \in \mathbb{R}$ and $x_i \in \mathcal{X}$ (positive definitness).
- $\langle K(x, \cdot), K(y, \cdot) \rangle_{\mathcal{H}} = K(x, y).$
- A Mercer kernel is a continuous positive definite kernel.

• The fact that Mercer kernels are positive definite and symmetric reminds us of similar properties of Gramians and covariance matrices. This is an essential fact that we are going to use in the following.

• Examples of kernels: $k(x, x') = \langle x, x' \rangle^d$, $k(x, x') = \exp\left(-\frac{||x-x'||}{2\sigma^2}\right)$, $k(x, x') = \tanh(\kappa \langle x, x' \rangle + \theta)$.

• Mercer Theorem: Let (\mathcal{X}, μ) be a finite-measure space, and suppose $k \in L_{\infty}(\mathcal{X}^2)$ is a symmetric real-valued function such that the integral operator

$$egin{array}{rl} T_k: L_2(\mathcal{X}) & o & L_2(\mathcal{X}) \ & f & \mapsto & (T_k f)(x) = \int_{\mathcal{X}} k(x,x') f(x') d\mu(x') \end{array}$$

is positive definite; that is, for all $f \in L_2(\mathcal{X})$, we have $\int_{\mathcal{X}^2} k(x,x') f(x) f(x') d\mu(x) d\mu(x') \geq 0$. Let $\Psi_j \in L_2(\mathcal{X})$ be the normalized orthogonal eigenfunctions of T_k associated with the eigenvalues $\lambda_j > 0$, sorted in non-increasing order. Then

i. $(\lambda_j)_j \in \ell_1$,

ii. $k(x, x') = \sum_{j=1}^{N_{\mathcal{X}}} \lambda_j \Psi_j(x) \Psi_j(x')$ holds for almost all (x, x'). Either $N_{\mathcal{X}} \in \mathbb{N}$, or $N_{\mathcal{X}} = \infty$; in the latter case, the series converges absolutely and uniformly for almost all $(x, x')_{\ominus}$.

• Proposition (Mercer Kernel Map): If k is a kernel satisfying the conditions in the preceding theorem, it is possible to construct a mapping Φ into a space where k acts as a dot product,

$$\langle \Phi(x), \Phi(x') \rangle = k(x, x'),$$

for almost all $x, x' \in \mathcal{X}$. Moreover, given any $\epsilon > 0$, there exists a map Φ_n into an n-dimensional dot product space (where $n \in \mathbb{N}$ depends on ϵ) such that

$$|k(x,x') - \langle \Phi(x), \Phi(x') \rangle| < \epsilon$$

for almost all $x, x' \in \mathcal{X}$.

• RKHS play an important in learning theory whose objective is to find an unknown function $f: X \to Y$ from random samples $(x_i, y_i)|_{i=1}^m$. • For instance, assume that the random probability measure that governs the random samples is ρ and is defined on $Z := X \times Y$. Let X be a compact subset of \mathbb{R}^n and $Y = \mathbb{R}$. If we define the least square error of f as $\mathcal{E} = \int_{X \times Y} (f(x) - y)^2 d\rho$, then the function that minimzes the error is the regression function $f_\rho f_\rho(x) = \int_{\mathbb{R}} y d\rho(y|x), \quad x \in X$, where $\rho(y|x)$ is the conditional probability measure on \mathbb{R} . • Since ρ is unknown, neither f_{ρ} nor \mathcal{E} is computable. We only have the samples $\mathbf{s} := (x_i, y_i)|_{i=1}^m$. The error f_ρ is approximated by the empirical error $\mathcal{E}_{s}(f)$ by

$$\mathcal{E}_{\mathbf{s}}(f) = \frac{1}{m} \sum_{i=1}^{m} (f(x_i) - y_i)^2 + \lambda ||f||_{\mathcal{H}}^2,$$

for $\lambda \geq 0$, λ plays the role of a regularizing parameter.

• In learning theory, the minimization is taken over functions from a hypothesis space often taken to be a ball of a RKHS \mathcal{H}_K associated to Mercer kernel K, and the function f_s that minimizes the empirical error \mathcal{E}_s is

$$f_{\mathbf{s}} = \sum_{j=1}^{m} c_j K(x, x_j),$$

where the coefficients $(c_j)_{j=1}^m$ is solved by the linear system

$$\lambda m c_i + \sum_{j=1}^m K(x_i, x_j) c_j = y_i, \quad i = 1, \cdots m,$$

and f_s is taken as an approximation of the regression function f_{ρ} . • We call *learning* the process of approximating the unknown function f from random samples on Z.

• We consider a general nonlinear system of the form

$$\left\{ \begin{array}{rrr} \dot{x} &=& f(x,u) \\ y &=& h(x) \end{array} \right.$$

with $x \in \mathbb{R}^n$, $u \in \mathbb{R}^m$, $y \in \mathbb{R}^p$, f(0,0) = 0, and h(0) = 0.

- Assume that the system is linear when lifted into an RKHS.
- In the linear case, $L_c(x_0) = \frac{1}{2}x_0^T W_c^{-1} x_0$ and $L_o(x_0) = \frac{1}{2}x_0^T W_o x_0$ can be rewritten as $L_c(x_0) = \frac{1}{2} \langle W_c^{\dagger} x_0, x_0 \rangle$ and $L_o(x_0) = \frac{1}{2} \langle W_o x_0, x_0 \rangle$. • In the nonlinear case, it may be tempting to write, in \mathcal{H} ,

 $L_c(h) = \frac{1}{2} \langle W_c^{\dagger}h, h \rangle$ and $L_o(h) = \frac{1}{2} \langle W_oh, h \rangle$. However, there are some complications...

• The domain of W_c^{\dagger} is equal to the range of W_c , and so in general K_x may not be in the domain of W_c^{\dagger} . We will therefore introduce the orthogonal projection $W_c^{\dagger}W_c$ mapping $\mathcal{H} \mapsto \operatorname{range}(W_c)$ and define the nonlinear control energy on \mathcal{H} as

$$L_c(h) = \left\langle W_c^{\dagger}(W_c^{\dagger}W_c)h, h \right\rangle.$$

• Since we will consider finite sample approximations to the preceding expression, $\widehat{W}_c^{\dagger}\widehat{W}_c$ may not converge to $W_c^{\dagger}W_c$ in the limit of infinite data (taking the pseudoinverse is not a continuous operation), and \widehat{W}_c^{\dagger} can easily be ill-conditioned in any event. One needs to impose regularization, and we replace the pseudoinverse A^{\dagger} with a regularized inverse $(A + \lambda I)^{-1}, \lambda > 0$ throughout.

• Intuitively, regularization prevents the estimator from overfitting to a bad or unrepresentative sample of data. We thus define the estimator $\hat{L}_c: \mathcal{X} \to \mathbb{R}_+$ (that is, on the domain $\{K_x \mid x \in \mathcal{X}\} \subseteq \mathcal{H}$) to be

$$\hat{L}_c(x) = \frac{1}{2} \big\langle (\widehat{W}_c + \lambda I)^{-2} \widehat{W}_c K_x, K_x \big\rangle, \quad x \in \mathcal{X}$$

with infinite-data limit

$$L_c^{\lambda}(x) = \frac{1}{2} \left\langle (W_c + \lambda I)^{-2} W_c K_x, K_x \right\rangle,$$

where $\lambda > 0$ is the regularization parameter.

- Towards deriving an equivalent but computable expression for \hat{L}_c defined in terms of kernels, we recall the sampling operator S_x introduced by Smale.
- Let $\mathbf{x} = \{x_i\}_{i=1}^m$ denote a generic sample of m data points. To \mathbf{x} we can associate the operators

$$S_{\mathbf{x}}: \mathcal{H} \to \mathbb{R}^m, \quad h \in \mathcal{H} \mapsto (h(x_1), \dots, h(x_m))$$
$$S_{\mathbf{x}}^*: \mathbb{R}^m \to \mathcal{H}, \qquad c \in \mathbb{R}^m \mapsto \sum_{i=1}^m c_i K_{x_i}.$$

If x is the collection of m = Np controllability samples, one can check that $\widehat{W}_c = \frac{1}{m} S^*_{\mathbf{x}} S_{\mathbf{x}}$ and $K_c = S_{\mathbf{x}} S^*_{\mathbf{x}}$.

• Consequently,

$$\hat{L}_{c}(x) = \frac{1}{2} \left\langle \left(\frac{1}{m} S_{\mathbf{x}}^{*} S_{\mathbf{x}} + \lambda I\right)^{-2} \frac{1}{m} S_{\mathbf{x}}^{*} S_{\mathbf{x}} K_{x}, K_{x} \right\rangle$$
$$= \frac{1}{2m} \left\langle S_{\mathbf{x}}^{*} \left(\frac{1}{m} S_{\mathbf{x}} S_{\mathbf{x}}^{*} + \lambda I\right)^{-2} S_{\mathbf{x}} K_{x}, K_{x} \right\rangle$$
$$= \frac{1}{2m} \mathbf{k}_{\mathbf{c}}(x)^{\top} \left(\frac{1}{m} K_{c} + \lambda I\right)^{-2} \mathbf{k}_{\mathbf{c}}(x),$$

where $\mathbf{k}_{\mathbf{c}}(x) := S_{\mathbf{x}}K_x = (K(x, x_{\mu}))_{\mu=1}^{Nq}$ is the Nq-dimensional column vector containing the kernel products between x and the controllability samples.

• Similarly, letting x now denote the collection of m = Np observability samples, we can approximate the future output energy by

$$\hat{L}_{o}(x) = \frac{1}{2} \langle \widehat{W}_{o} K_{x}, K_{x} \rangle$$

$$= \frac{1}{2m} \langle S_{\mathbf{x}}^{*} S_{\mathbf{x}} K_{x}, K_{x} \rangle$$

$$= \frac{1}{2m} \mathbf{k}_{o}(x)^{\mathsf{T}} \mathbf{k}_{o}(x) = \frac{1}{2m} \| \mathbf{k}_{o}(x) \|_{2}^{2}$$
(1)

where $\mathbf{k}_{\mathbf{o}}(x) := (K(x, d_{\mu}))_{\mu=1}^{Np}$ is the Np-dimensional column vector containing the kernel products between x and the observability samples.

Balanced Reduction of Nonlinear Control Systems in RKHS

• We consider a general nonlinear system of the form

$$\begin{cases} \dot{x} &= f(x, u) \\ y &= h(x) \end{cases}$$

with $x \in \mathbb{R}^n$, $u \in \mathbb{R}^m$, $y \in \mathbb{R}^p$, f(0,0) = 0, and h(0) = 0. We assume that the system is zero-state observable, and that the linearization of around the origin is controllable. We also assume that the origin of $\dot{x} = f(x,0)$ is asymptotically stable.

Proposed Data-Driven Approach:

- Assume that the system behaves linearly when lifted to a high dimensional feature space.
- Carry out balancing and truncation (linear techniques) implicitly in the feature space (discard unimportant states).

 Idea: We can perform balancing/truncation in feature space by lifting the data into \mathcal{H} via Φ , and simultaneously diagonalizing the corresponding covariance operators.

The empirical controllability Gramian

$$\widehat{W}_{c} = \frac{T}{mN} \sum_{i=1}^{N} X(t_{i}) X(t_{i})^{\top} = \frac{T}{mN} \sum_{i=1}^{N} \sum_{j=1}^{m} x^{j}(t_{i}) x^{j}(t_{i})^{\top}$$

becomes

$$C_c = \frac{T}{mN} \sum_{i=1}^{N} \sum_{j=1}^{m} \left\langle \Phi\left(x^j(t_i)\right), \cdot \right\rangle_{\mathcal{H}} \Phi\left(x^j(t_i)\right)$$

for example.

Balancing in RKHS

• "Balancing" is carried out implicitly in \mathcal{H} by simultaneous diagonalization of K_c and K_o .

• If $K_c^{1/2}K_oK_c^{1/2} = U\Sigma^2 U^{\mathsf{T}}$, we can define the aligning transformation

$$T = \Sigma^{1/2} U^{\mathsf{T}} \sqrt{K_c^{\dagger}}.$$

• The dimension of the state space is reduced by discarding small eigenvalues $\{\Sigma_{ii}\}_{i=q+1}^{n}$, and projecting onto the subspace in \mathcal{H} associated with the first q < n largest eigenvalues.

• This leads to the *nonlinear* state-space dimensionality reduction map $\Pi:\mathbb{R}^n\to\mathbb{R}^q$ given by

$$\Pi(x) = T_q^{\mathsf{T}} \mathbf{k}_c(x), \quad x \in \mathbb{R}^n$$

where

$$\mathbf{k}_c(x) := \left(K(x, x^1(t_1)), \dots, K(x, x^m(t_N)) \right)^{\top}.$$

Consider the 7 - D system (Nilsson, 2009)

$$\begin{aligned} \dot{x}_1 &= -x_1^3 + u & \dot{x}_2 &= -x_2^3 - x_1^2 x_2 + 3x_1 x_2^2 - u \\ \dot{x}_3 &= -x_3^3 + x_5 + u & \dot{x}_4 &= -x_4^3 + x_1 - x_2 + x_3 + 2u \\ \dot{x}_5 &= x_1 x_2 x_3 - x_5^3 + u & \dot{x}_6 &= x_5 - x_6^3 - x_5^3 + 2u \\ \dot{x}_7 &= -2x_6^3 + 2x_5 - x_7 - x_5^3 + 4u \\ y &= x_1 - x_2^2 + x_3 + x_4 x_3 + x_5 - 2x_6 + 2x_7 \end{aligned}$$

э

- Excite with impulses: inputs (K_c) and initial conditions $(K_o, u = 0)$.
- ► Learn \hat{f}, \hat{h} using a 10Hz square wave input signal u.
- Reduce to a second-order system.
- Simulate the reduced system with a different input,

$$u(t) = \frac{1}{2} \left(\sin(2\pi 3t) + \mathsf{sq}(2\pi 5t - \pi/2) \right)$$

and compare the output to that of the original system.

Experiment



э June 4th, 2012 38 / 55

< 同 ▶

э

Experiment



< 1 →

æ

- Consider the stochastically excited stable dynamical control systems affine in the input $u\in \mathbb{R}^q$

$$\dot{x} = f(x) + G(x)u \; ,$$

where $G : \mathbb{R}^n \to \mathbb{R}^{n \times q}$ is a smooth matrix-valued function. We replace the control inputs by sample paths of white Gaussian noise processes, giving the corresponding stochastic differential equation (SDE)

$$dX_t = f(X_t)dt + G(X_t)dW_t^{(q)}$$

with $W_t^{(q)}$ a q-dimensional Brownian motion. The solution X_t to this SDE is a Markov stochastic process with transition probability density $\rho(t, x)$ that satisfies the Fokker-Planck (or Forward Kolmogorov) equation

$$\frac{\partial \rho}{\partial t} = -\langle \frac{\partial}{\partial x}, f\rho \rangle + \frac{1}{2} \sum_{j,k=1}^{n} \frac{\partial^2}{\partial x_j \partial x_k} [(GG^T)_{jk}\rho] =: L\rho \; .$$

- In the context of linear Gaussian theory where we are given an n-dimensional system of the form $dX_t = AX_t dt + BdW_t^{(q)}$, with $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times q}$, the transition density is Gaussian.
- It is therefore sufficient to find the mean and covariance of the solution X(t) in order to uniquely determine the transition probability density.

- The mean satisfies $\frac{d}{dt}\mathbb{E}[x] = A\mathbb{E}[x]$ and thus $\mathbb{E}[x(t)] = e^{At}\mathbb{E}[x(0)]$. If A is Hurwitz, $\lim_{t\to\infty}\mathbb{E}[x(t)] = 0$.
- The covariance satisfies $\frac{d}{dt}\mathbb{E}[xx^T] = A\mathbb{E}[xx^T] + \mathbb{E}[xx^T]A + BB^T$.
- Hence, $\mathcal{Q} = \lim_{t \to \infty} \mathbb{E}[xx^{\top}]$ satisfies the Lyapunov system $A\mathcal{Q} + \mathcal{Q}A^{\top} = -BB^{\top}$. So, $\mathcal{Q} = W_c = \int_0^\infty e^{At}BB^{\top}e^{A^{\top}t} dt$, where W_c is the controllability Gramian, which is positive iff. the pair (A, B) is controllable.

- Combining the above facts, the steady-state probability density is given by $\rho_\infty(x)=Z^{-1}e^{-\frac12x^\top W_c^{-1}x}=Z^{-1}e^{-L_c(x)}$

with $Z = \sqrt{(2\pi)^n \det(W_c)}$.

• The preceding suggests the following key observations in the linear setting: Given an approximation \hat{L}_c of L_c we obtain an approximation for ρ_∞ of the form

 $\hat{\rho}_{\infty}(x) \propto e^{-\hat{L}_c(x)}$

• Although the above relationship between ρ_{∞} and L_c holds for only a small class of systems (e.g. linear and some Hamiltonian systems), by mapping a nonlinear system into a suitable reproducing kernel Hilbert space we may reasonably extend this connection to a broad class of nonlinear systems.

• Assumption1: Given a suitable choice of kernel K, if the \mathbb{R}^d -valued stochastic process x(t) is a solution to the (ergodic) stochastically excited nonlinear system

$$dX_t = f(X_t)dt + G(X_t)dW_t^{(q)}$$

the $\mathcal H\text{-valued}$ stochastic process $(\Phi\circ x)(t)=:X(t)$ can be reasonably modelled as an Ornstein-Uhlenbeck process

$$dX(t) = AX(t)dt + \sqrt{C}dW(t), \quad X(0) = 0 \in \mathcal{H}$$

where A is linear, negative and is the infinitesimal generator of a strongly continuous semigroup e^{tA} , C is linear, continuous, positive and self-adjoint, and W(t) is the cylindrical Wiener process.

• Assumption2: The measure P_{∞} is the invariant measure of the OU process and P_{∞} is the pushforward along Φ of the unknown invariant measure μ_{∞} on the statespace \mathcal{X} we would like to approximate. • Assumption3: The measure μ_{∞} is absolutely continuous with respect to Lebesgue measure, and so admits a density. • The stationary measure μ_∞ is defined on a finite dimensional space, so together with part (iii) of Assumption A, we may consider the corresponding density

 $\rho_{\infty}(x) \propto \exp\left(-\hat{L}_c(x)\right)$

Experiment

Consider the SDE $dX = -5X^5 + 10X^3 + \sqrt{2}dW$.



- It is apparent from $L\rho_{\infty}(x) = 0$ that both drift and diffusion coefficients of an SDE are directly related to the stationary solution ρ_{∞} of the Fokker-Planck equation.
- This relation can be employed to derive estimators for the unknown coefficients.

Parameter Estimation for SDEs

• Consider the scalar-valued SDE

$$dX_t = f(X_t) dt + g(X_t) dW_t .$$

• Equation $L\rho_{\infty}(x) = 0$ reduces to

$$f(x)\rho_{\infty}(x) = \frac{1}{2}(g(x)^2\rho_{\infty}(x))'$$
.

The methodology described before yields an estimator $\hat{\rho}_{\infty}$ of ρ_{∞} , but this is obviously not sufficient to estimate both unknown functions f and g directly using the above relation. If we, however, have knowledge of either one, then a natural nonparametric estimator for the other one is given via

$$\hat{f}(x) = g(x)g'(x) + \frac{g(x)^2\hat{\rho}'_{\infty}(x)}{2\hat{\rho}_{\infty}(x)}$$

or

$$\hat{g}(x)^2 = \frac{2}{\hat{\rho}_{\infty}(x)} \int_0^x f(u)\hat{\rho}_{\infty}(u) \, du \; ,$$

• It is reasonable to assume that the diffusion coefficient is known up to a multiplicative factor $g(x; \theta) = \theta g(x)$ with θ being an unknown parameter and g is known here. In this case we estimate θ via the quadratic variation of the path estimator $\hat{\theta}$

$$\hat{\theta} = \frac{\sum_{i=0}^{n-1} (X_{t_{i+1}} - X_{t_i})^2}{\int_0^t g(X_s) \, ds}$$

Example (Fokker-Planck Integrable Systems)

Consider the SDE

$$dx_t = (\alpha x_t + \beta x_t^3) \, dt + \sqrt{\sigma} \, dW_t$$

• Provided α and β are such that $\lim_{|x|\to\infty} \Phi(x) = \infty$ and $e^{-\gamma\Phi} \in L^1(\mathbb{R})$ for all $\gamma > 0$, the invariant density is given by

$$\rho_{\infty}(x) = Z e^{-\frac{2}{\sigma}\Phi(x)} = Z e^{\frac{2}{\sigma}(\frac{\alpha}{2}x^2 + \frac{\beta}{4}x^4)} ,$$

with Z being the appropriate normalization constant.

Example (A Fokker-Planck Integrable System)

• The stationary Fokker-Planck equation for this example reads $f(x) = \sigma \rho'_{\infty}(x)/(2\rho_{\infty}(x))$, so that we find $\alpha x + \beta x^3 = \frac{\sigma}{2} \frac{\rho'_{\infty}(x)}{\rho_{\infty}(x)}, \quad \forall x \in \operatorname{supp}(\rho_{\infty}).$

• To obtain the parameters as a least squares fit to the observations, let $\{x_i\}_{i=1}^m \subset \mathbb{R}$, with $x_i \in \operatorname{supp}(\hat{\rho}_{\infty})$, denote a (finite) sequence of samples. Since the preceding equation holds for all $x \in \operatorname{supp}(\rho_{\infty}) \supset \operatorname{supp}(\hat{\rho}_{\infty})$ we have that

$$\alpha x_i + \beta x_i^3 = \frac{\hat{\sigma}}{2} \frac{\hat{\rho}'_{\infty}(x_i)}{\hat{\rho}_{\infty}(x_i)} ,$$

for $1 \le i \le m$. Consequently, the estimators $\hat{\alpha}$ and $\hat{\beta}$ obtained by a least squares fit solve the system of linear equations

$$\begin{pmatrix} \sum_{i=1}^{m} x_i^2 & \sum_{i=1}^{m} x_i^4 \\ \sum_{i=1}^{m} x_i^4 & \sum_{i=1}^{m} x_i^6 \end{pmatrix} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \frac{\hat{\sigma}}{2} \sum_{i=1}^{m} \frac{\hat{\rho}'_{\infty}(x_i)}{\hat{\rho}_{\infty}(x_i)} \begin{pmatrix} x_i \\ x_i^3 \end{pmatrix}$$

which can be solved explicitly provided the system_matrix is_invertible.

- We have introduced estimators for the controllability/observability energies and reachable/observable sets of nonlinear control systems.
- We showed that the controllability energy estimator may be used to estimate the stationary solution (and its support) of the Fokker-Planck equation governing nonlinear SDEs.
- The estimators we derived were based on applying linear methods for control and random dynamical systems to nonlinear control systems and SDEs, once mapped into an infinite-dimensional RKHS acting as a "linearizing space".
- These results collectively argue that there is a reasonable passage from linear dynamical systems theory to a data-based nonlinear dynamical systems theory through reproducing kernel Hilbert spaces.

- Derivation of data-based estimators for Lyapunov exponents and the controllability and the observability operators.
- Once a data-based approximation of the controllability operator is obtained, derive a data-based controller.
- We have been using Mercer kernels for our experiments: is there a better way to find good kernels ? (kernel choice as an optimization problem: kernel as encoding "data + some information on the dynamics (constraint)")
- Explicitly compute the embedings $\Phi : \mathbb{R}^n \mapsto \mathcal{H}$ and writing down the dynamics in \mathcal{H} .
- Error estimates when more knowledge about the dynamics is given.
- How to use methods from Linear Systems Identification to perform parameter estimation of nonlinear systems in RKHSes.