

EM Algorithms for Learning Local Dynamics Models in a Bayesian Framework

Samuel Otto

Outline

1. Motivation
2. Problem Setup
3. Maximum Likelihood and the EM Algorithm
4. Explicit Solutions for Maximization Step
 - i. Weighted Regression
 - ii. Weighted Gaussian Mixture
 - iii. Categorical “Multinoulli” Distribution
5. Reconstructing the Estimates
6. Testing on a Chaotic Poincare Map
7. K-Means and Generalized K-Means as efficient EM variants
 - i. Test on Chaotic Poincare Map
8. Extending to Nonlinear Local Dynamics using Kernels
 - i. Test Variants on Chaotic Poincare Map
9. Extending to Non-Gaussian Density Estimation using Kernels
10. Future Work

Motivation

- Model complex, possibly chaotic dynamics using data.
 - Effect of control perturbations
 - Short-term prediction
 - Generative modeling: plausible dynamics sharing statistical properties
- Analyze dynamics locally and globally
 - Local spectral analysis
 - Identify global structures like orbits, tori, and attractors
- Understand the distribution of data points in phase space
 - Structure of data distribution in phase space
 - Identify regions of interest

Three Main Approaches

- If dynamics vary nonlinearly over phase space, we have three main approaches for modeling:
 1. Use a single model with many nonlinear terms giving it enough representational capacity to capture dynamics everywhere.
 - EDMD: Use large number of nonlinear observables
 - Kernel Method: Kernel matrix grows with the square of the number of data points
 2. Piece together many local models with fewer nonlinear terms in each
 - Models with fewer parameters can be updated and evaluated quickly
 - Representational power is increased simply by adding more models
 3. Fully non-parametric representation using locally-weighted regression (Not practical)

Local Modeling Approach

- Model distribution of data in phase space using mixture of Gaussians
 - Universal approximator for continuous distributions
- To each Gaussian, we associate a model of the dynamics in that neighborhood
- Introduce a latent variable indicating which model describes the dynamics at a given point.
- Allows for new points to be classified according to maximum likelihood among the models
 - Finds the model most likely to capture dynamics at the new point
- Must learn the mixture of Gaussians and corresponding local dynamics together
 - Data clustering alone does not reflect the number of local models needed to capture dynamics in a region

Problem Setup

- Suppose we wish to approximate $y = f(x, u)$, $x, y \in \mathbb{R}^n$, $u \in \mathbb{R}^q$ with training data $\{(x_j, u_j, y_j)\}_{j=1}^m$
- Introduce N models $\mathbf{M} = \{M^i = (f^i, \mu_x^i, \mu_y^i, R^i, \Sigma^i, \phi^i)\}_{i=1}^N$ of the form

$$Y^i = f^i(x, u) + V^i, \quad V^i \sim p_{V^i}(v^i) = \mathcal{N}(v^i, \underline{0}, R^i)$$

- Introduce latent random variable $Z \in \{1, \dots, N\}$ indicating the model. The prior probabilities are assumed to be “Multinoulli”

$$P(Z = i) = \phi^i$$

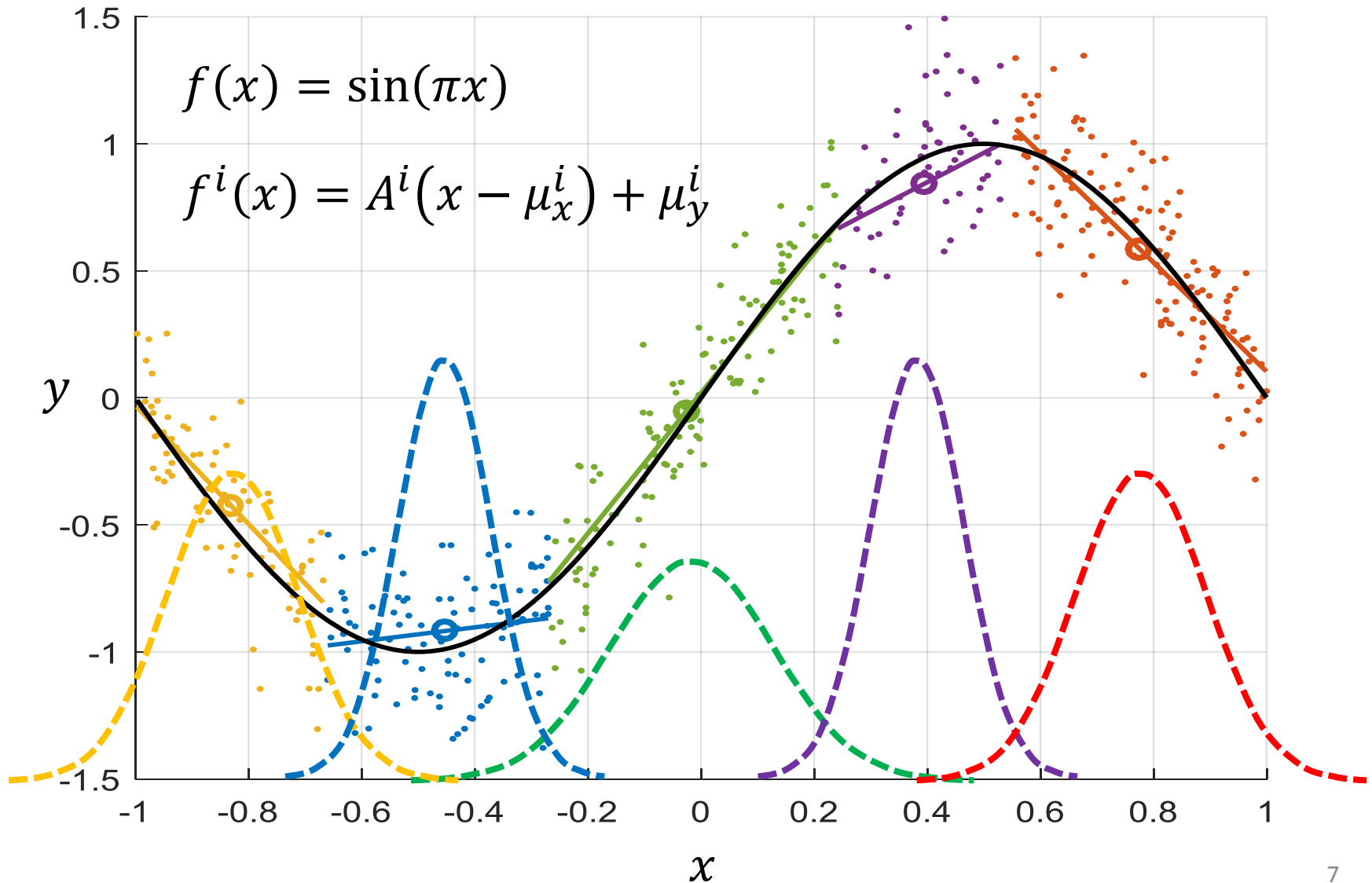
- Define posterior probabilities of the Gaussian mixture:

$$p_{X|Z=i}(x|Z = i) = \mathcal{N}(x, \mu_x^i, \Sigma^i)$$

- Infer conditional model probabilities using Bayes' Rule

$$P(Z = i|X = x) = \frac{p_{X|Z=i}(x|Z = i)P(Z = i)}{\sum_{k=1}^N p_{X|Z=k}(x|Z = k)P(Z = k)}$$

Problem Setup



Maximum Likelihood Estimation

- Likelihood function based on joint probability (m is number of training points)

$$L(\mathbf{M}) = \prod_{j=1}^m P_{model}(y_j, x_j; \mathbf{M})$$

- Training on log likelihood is equivalent to minimizing cross entropy

$$l(\mathbf{M}) = \mathbb{E}_{X,Y \sim \hat{P}_{data}} [\log P_{model}(X, Y; \mathbf{M})] = \frac{1}{m} \sum_{j=1}^m \log P_{model}(y_j, x_j; \mathbf{M})$$

- Condition on latent variable Z and use total probability (N is number of models)

$$P_{model}(y_j, x_j; \mathbf{M}) = \sum_{i=1}^N P(y_j, x_j, Z = i)$$

- Chain rule for conditional probability

$$P_{model}(y_j, x_j; \mathbf{M}) = \sum_{i=1}^N p_{Y|Z,X,u}(y_j | u_j, x_j, Z = i) p_{X|Z=i}(x_j | Z = i) P(Z = i)$$

$$P_{model}(y_j, x_j; \mathbf{M}) = \sum_{i=1}^N p_{V^i}(y_j - f^i(x_j, u_j)) p_{X|Z=i}(x_j | Z = i) P(Z = i)$$

Maximum Likelihood Estimation

- Introduce a yet to be determined distribution over the models at each training point $Q_j(Z)$

$$l(\mathbf{M}) = \frac{1}{m} \sum_{j=1}^m \log \left[\sum_{i=1}^N Q_j(Z = i) \frac{p_{V^i}(v_j^i) p_{X|Z=i}(x_j|Z = i) P(Z = i)}{Q_j(Z = i)} \right]$$

- Observe that the inner sum is an expectation

$$l(\mathbf{M}) = \frac{1}{m} \sum_{j=1}^m \log \left(\mathbb{E}_{Z \sim Q_j} \left[\frac{p_{V^Z}(v_j^Z) p_{X|Z}(x_j|Z) P(Z)}{Q_j(Z)} \right] \right)$$

- Use Jensen's inequality to lower bound the log likelihood

$$l(\mathbf{M}) \geq \tilde{l}(\mathbf{M}) = \frac{1}{m} \sum_{j=1}^m \mathbb{E}_{Z \sim Q_j} \left[\log \left(\frac{p_{V^Z}(v_j^Z) p_{X|Z}(x_j|Z) P(Z)}{Q_j(Z)} \right) \right]$$

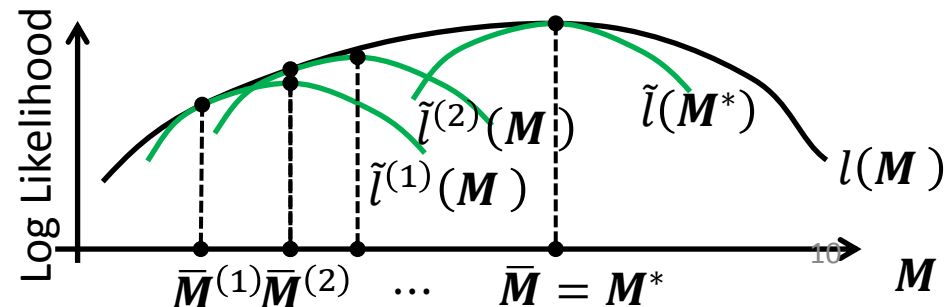
$$l(\mathbf{M}) \geq \tilde{l}(\mathbf{M}) = \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^N Q_j(Z = i) \log \left(\frac{p_{V^i}(v_j^i) p_{X|Z=i}(x_j|Z = i) P(Z = i)}{Q_j(Z = i)} \right)$$

- Important note: this inequality becomes an equality when the inside term does not depend on Z

The EM Algorithm

- Iterative update process for the model parameters \mathbf{M} given lower bounds using previous set of parameters $\bar{\mathbf{M}}$.
 - Maximizing lower bounds $\tilde{l}(\mathbf{M})$ gives monotone increase in log likelihood function.
- Clever choice for $Q_j(Z)$ such that the inequality becomes equality when the algorithm has converged $\bar{\mathbf{M}} = \mathbf{M}^*$.

$$Q_j(Z = i; \bar{\mathbf{M}}) = \frac{p_{Vi}(v_j^i; \bar{\mathbf{M}})p_{X|Z=i}(x_j|Z = i; \bar{\mathbf{M}})P(Z = i; \bar{\mathbf{M}})}{\sum_{k=1}^N p_{Vk}(v_j^k; \bar{\mathbf{M}})p_{X|Z=k}(x_j|Z = k; \bar{\mathbf{M}})P(Z = k; \bar{\mathbf{M}})}$$



The EM Algorithm

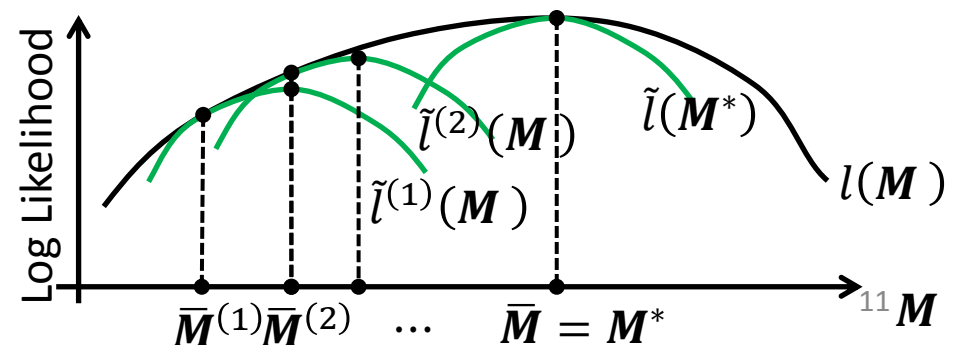
1. Expectation Step: Given previous set of parameter estimates $\bar{\mathbf{M}}$, compute the probabilities $Q_j(Z = i; \bar{\mathbf{M}})$ to use in the expectation

$$\bar{Q}_j^i \triangleq Q_j(Z = i; \bar{\mathbf{M}}) = \frac{p_{Vi}(v_j^i; \bar{\mathbf{M}})p_{X|Z=i}(x_j|Z = i; \bar{\mathbf{M}})P(Z = i; \bar{\mathbf{M}})}{\sum_{k=1}^N p_{Vk}(v_j^k; \bar{\mathbf{M}})p_{X|Z=k}(x_j|Z = k; \bar{\mathbf{M}})P(Z = k; \bar{\mathbf{M}})}$$

2. Maximization Step: Maximize the lower bound $\tilde{l}(\mathbf{M})$ using the above probabilities in the expectation

$$\bar{\mathbf{M}} \leftarrow \operatorname{argmax}_{\mathbf{M}} \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \log \left(\frac{p_{Vi}(v_j^i; \mathbf{M})p_{X|Z=i}(x_j|Z = i; \mathbf{M})P(Z = i; \mathbf{M})}{\bar{Q}_j^i} \right)$$

3. Iterate until convergence



Maximization Step

- Helpful because we now have the log of a product which decouples into separate optimizations. Substituting the known forms we have

$$\bar{\mathbf{M}} \leftarrow \operatorname{argmax}_{\mathbf{M}} \underbrace{\sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \log[\mathcal{N}(y_j, f^i(x_j, u_j), R^i)]}_{(I)} + \underbrace{\sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \log[\mathcal{N}(x_j, \mu_x^i, \Sigma^i)]}_{(II)} + \underbrace{\sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \log(\phi^i)}_{(III)}$$

I. Weighted regression

$$\operatorname{argmax}_{\mathbf{M}}(I) = \operatorname{argmax}_{\{(f^i, R^i)\}_{i=1}^N} \sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \left[-\frac{1}{2} \log |R^i| - \frac{1}{2} (y_j - f^i(x_j, u_j))^T (R^i)^{-1} (y_j - f^i(x_j, u_j)) \right]$$

II. Weighted mixture of Gaussians

$$\operatorname{argmax}_{\mathbf{M}}(II) = \operatorname{argmax}_{\{(\mu_x^i, \Sigma^i)\}_{i=1}^N} \sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \left[-\frac{1}{2} \log |\Sigma^i| - \frac{1}{2} (x_j - \mu_x^i)^T (\Sigma^i)^{-1} (x_j - \mu_x^i) \right]$$

III. Maximum likelihood “Multinoulli”

$$\operatorname{argmax}_{\mathbf{M}}(III) = \operatorname{argmax}_{\{\phi^i\}_{i=1}^N} \sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \log \phi^i \quad s.t. \quad 1 = \sum_{i=1}^N \phi^i$$

(I) Weighted Regression

- Form vector with (possibly nonlinear) observables $\theta_1, \theta_2, \dots, \theta_L: \mathbb{R}^{n+q} \rightarrow \mathbb{R}$ of the state and control

$$\Phi_j \triangleq [\theta_1(x_j, u_j) \quad \dots \quad \theta_L(x_j, u_j)]^T \quad \text{or} \quad \Phi_j \triangleq [x_j \quad u_j \quad 1]^T$$

- The learned functions are then linear combinations of observables

$$f^i(x, u) = A^i \Phi(x, u) = A^i [\theta_1(x_j, u_j) \quad \dots \quad \theta_L(x_j, u_j)]^T$$

- Take derivatives and do math to optimize (I) analytically.

$$A^i = \left(\sum_{j=1}^m \bar{Q}_j^i y_j \Phi_j^T \right) \left(\sum_{j=1}^m \bar{Q}_j^i \Phi_j \Phi_j^T \right)^+ \quad R^i = \frac{\sum_{j=1}^m \bar{Q}_j^i (y_j - A^i \Phi_j) (y_j - A^i \Phi_j)^T}{\sum_{j=1}^m \bar{Q}_j^i}$$

- Letting $\mathbf{Y} = [y_1 \dots y_m]$, $\mathbf{\Phi} = [\Phi_1 \dots \Phi_m]$, and $W^i = \text{diag}[\bar{Q}_1^i, \dots, \bar{Q}_m^i]$ we have the familiar weighted least squares equations

$$A^i = (\mathbf{Y} W^i \mathbf{\Phi}^T) (\mathbf{\Phi} W^i \mathbf{\Phi}^T)^+ \quad R^i = (\mathbf{1}_m^T W^i \mathbf{1}_m)^{-1} (\mathbf{Y} - A^i \mathbf{\Phi}) W^i (\mathbf{Y} - A^i \mathbf{\Phi})^T$$

(II) Weighted Mixture of Gaussians

- Taking derivatives and optimizing (II) gives the weighted mixture of Gaussians model
- Let $\mathbf{X} = [x_1 \cdots x_m]$ and $W^i = \text{diag}[\bar{Q}_1^i, \dots, \bar{Q}_m^i]$
- The centers are found:

$$\mu_x^i = \frac{\sum_{j=1}^m \bar{Q}_j^i x_j}{\sum_{j=1}^m \bar{Q}_j^i} = \frac{\mathbf{X} W^i \mathbf{1}_m}{\mathbf{1}_m^T W^i \mathbf{1}_m}$$

- The covariances are found:

$$\Sigma^i = \frac{\sum_{j=1}^m \bar{Q}_j^i (x_j - \mu_x^i)(x_j - \mu_x^i)^T}{\sum_{j=1}^m \bar{Q}_j^i} = \frac{(\mathbf{X} - \mu_x^i \mathbf{1}_m^T) W^i (\mathbf{X} - \mu_x^i \mathbf{1}_m^T)^T}{\mathbf{1}_m^T W^i \mathbf{1}_m}$$

(III) Weighted Multinoulli

- Form the Lagrangian to include the constraint

$$\mathcal{L}(\phi^1, \dots, \phi^N, \lambda) = \sum_{j=1}^m \sum_{i=1}^N \bar{Q}_j^i \log \phi^i + \lambda \left(1 - \sum_{i=1}^N \phi^i \right)$$

- Setting derivatives to zero and doing some algebra gives the solution

$$\boxed{\phi^i = \frac{1}{m} \sum_{j=1}^m \bar{Q}_j^i}$$

- In the expressions for the other solutions, we have the normalization constant $\mathbf{1}_m^T W^i \mathbf{1}_m = \sum_{j=1}^m \bar{Q}_j^i = m\phi^i$ in the denominator
 - This is the expected number of training points belonging to each model a priori.

Reconstructing the Dynamics

- Given a new point x and control input u , how do we predict the dynamics?
- Conditional model probabilities using Bayes' rule

$$P(Z = i|X = x) = \frac{\mathcal{N}(x, \mu_x^i, \Sigma^i) \phi^i}{\sum_{k=1}^N \mathcal{N}(x, \mu_x^k, \Sigma^k) \phi^k}$$

1. Maximum likelihood model assignment

$$\hat{y} = f^{i^*}(x, u) \quad i^* = \underset{i}{\operatorname{argmax}} P(Z = i|X = x)$$

2. Stochastic model assignment

$$\hat{y} = f^{i^*}(x, u) \quad i^* \sim P(Z|X = x)$$

3. Weighted model assignment

$$\hat{y} = \mathbb{E}[Y] = \sum_{i=1}^N P(Z = i|X = x) \mathbb{E}[Y^i] = \sum_{i=1}^N P(Z = i|X = x) f^i(x, u)$$

4. Hybrid weighted model assignment from top $N' \leq N$ models

- Sort models (i_1, i_2, \dots, i_N) from greatest to least probability $P(Z|X = x)$

$$\hat{y} = \left(\sum_{l=1}^{N'} P(Z = i_l|X = x) \right)^{-1} \sum_{l=1}^{N'} P(Z = i_l|X = x) f^{i_l}(x, u)$$

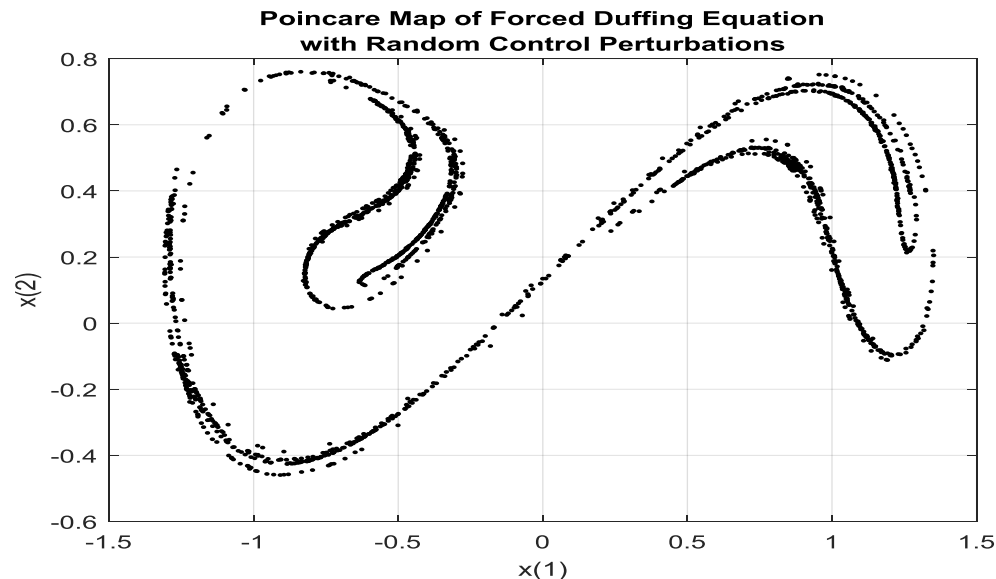
5. Some kind of Kriging may also be possible. This is a topic for future work.

Application to Duffing Equation with Control

- Forced Duffing equation with control input

$$\frac{d}{dt} \begin{bmatrix} x(1) \\ x(2) \end{bmatrix} = \begin{bmatrix} x(2) \\ \gamma \cos(\omega t) - \delta x(2) - \alpha x(1) - \beta x(1)^3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 1 & \cos \omega t & \sin \omega t & \cos 2\omega t & \sin 2\omega t \end{bmatrix} \begin{bmatrix} u(1) \\ u(2) \\ u(3) \\ u(4) \\ u(5) \end{bmatrix}$$

- Define time $T = 2\pi/\omega$ Poincare map with constant control u over each interval $x_{n+1} = f(x_n, u_n)$
- The Poincare map was simulated for $m = 2000$ iterations with small uniformly distributed random control perturbations $u(k) \in [-10^{-3}, 10^{-3}]$



Error Metrics

- The algorithms will be evaluated on $m_{test} = 1000$ new testing points generated from new random control perturbations
- Mean estimation error distance

$$\bar{E}_D = \frac{1}{m_{test}} \sum_{j=1}^{m_{test}} \|y_j - \hat{y}_j\|$$

- Decaying error norm after many steps. Start at x_0 and iterate the map f and approximate map \hat{f} with random control inputs u_0, u_2, \dots, u_{M-1} generating sequences x_1, x_2, \dots, x_M and $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_M$ respectively

$$E_\alpha^M(x_0) = \sum_{n=1}^M \frac{\|x_n - \hat{x}_n\|}{\alpha^{n-1}}$$

- The average is taken over the data distribution by sampling many initial points

$$\bar{E}_\alpha^M = \frac{1}{m_{test}} \sum_{j=1}^{m_{test}} E_\alpha(x_j)$$

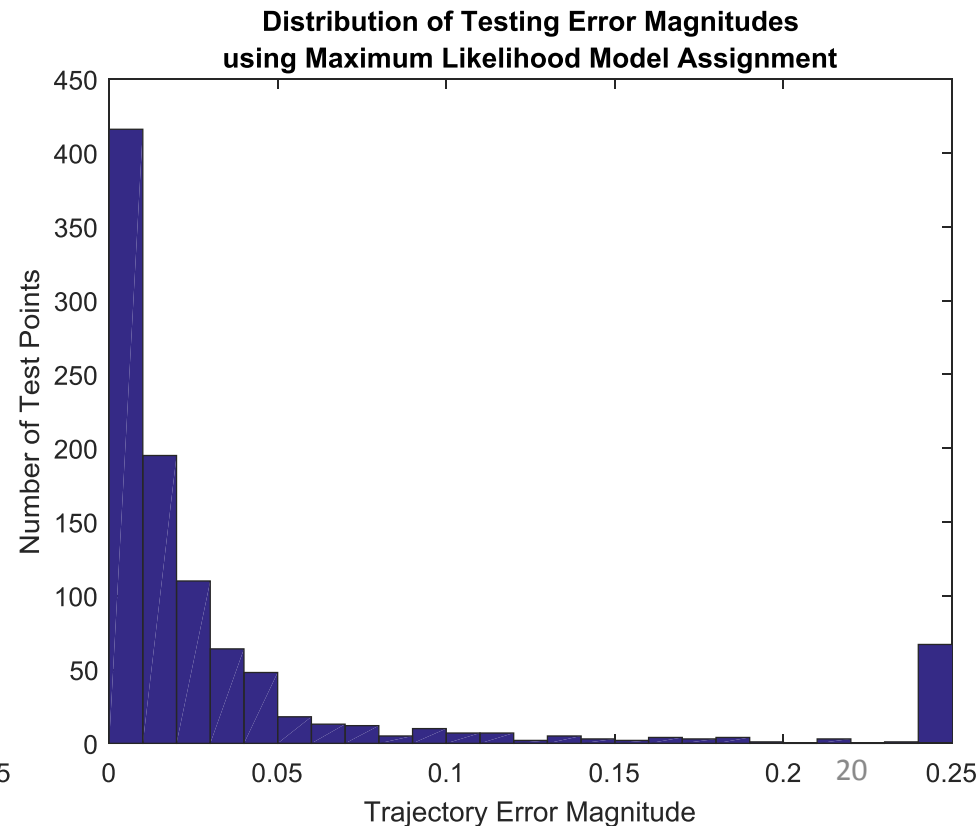
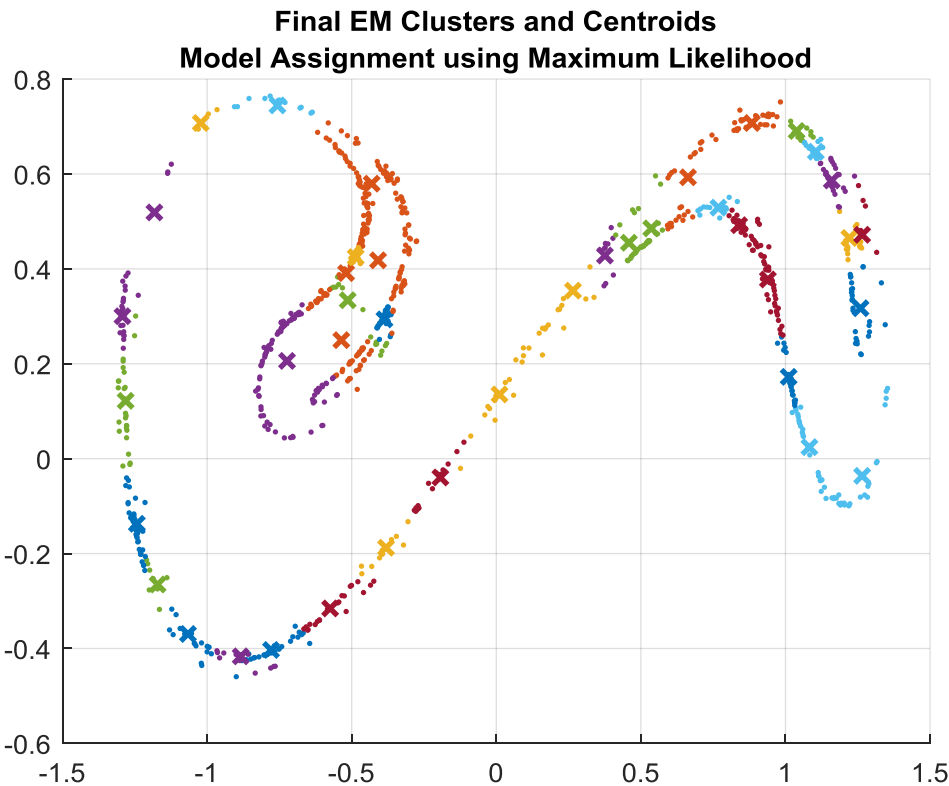
- The autocorrelations of the uncontrolled dynamics over many iterations will be compared in order to estimate stationary statistics $R(n) = \mathbb{E}[x_k x_{k-n}^T]$ of the actual and $\hat{R}(n) = \mathbb{E}[\hat{x}_k \hat{x}_{k-n}^T]$ of the modeled dynamics.
- The histograms of the distance errors on the 1000 test points will also be considered

Linear EM Gaussian Mixture Results

- $N = 40$ linear models were fit
- Each model is fit to the $m' = 500$ training points with the largest weights \bar{Q}_j^i
- Maximum likelihood model assignment was used to reconstruct the dynamics
 - Numerical experiments show that of the methods discussed, maximum likelihood produced the lowest errors.
- Minimum singular values of covariances Σ and R were needed to prevent degeneracy
- $m_{test} = 1000$ independently generated test points with random control perturbations were used for evaluation

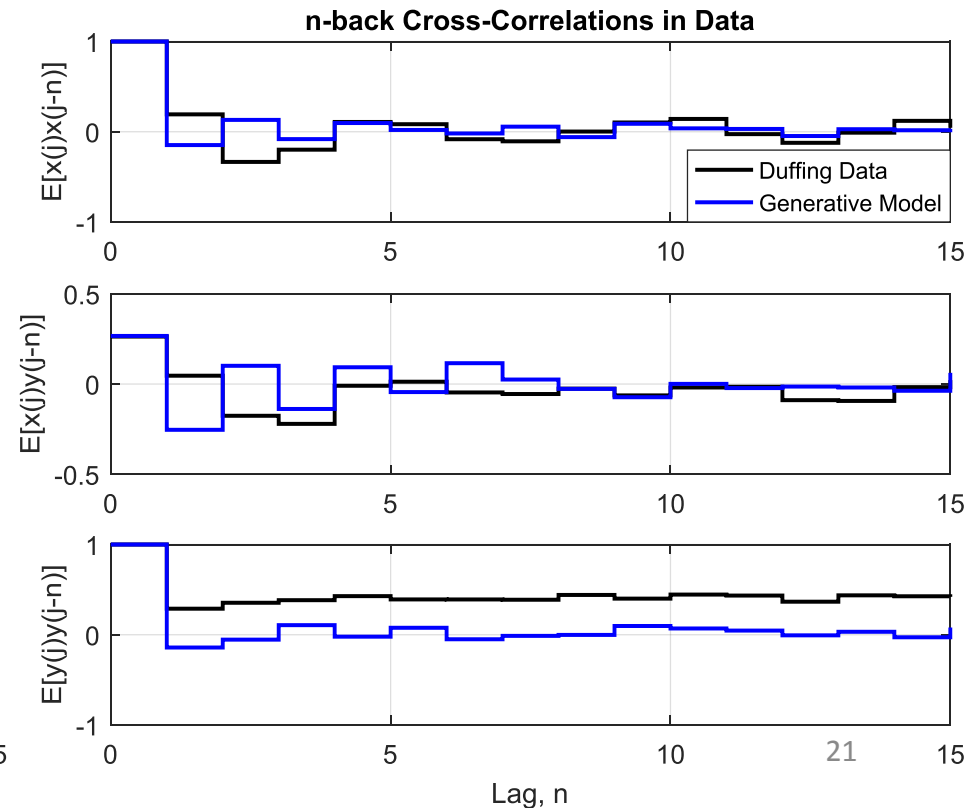
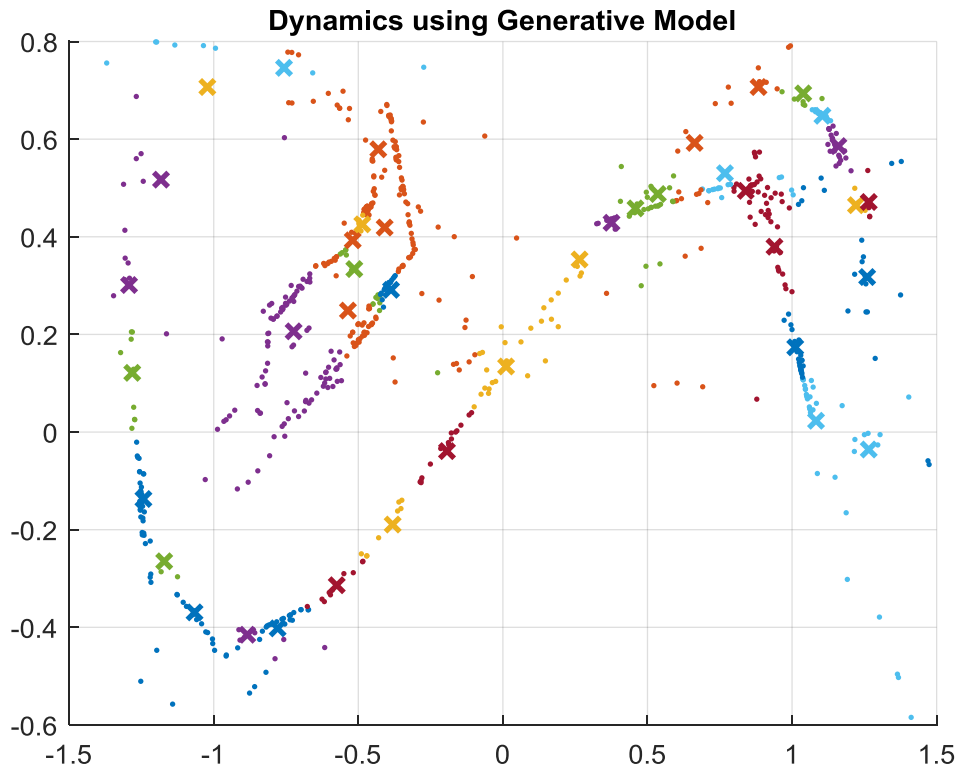
Linear EM Gaussian Mixture Results

- Mean Distance Error: $\bar{E}_D = 0.0608$
- Median Distance Error: 0.0134
- Decaying Error Norm over $M = 10$ steps with $\alpha = 2$: $\bar{E}_\alpha^M = 0.3921$



Use as a Generative Model

- The dynamics were simulated using the learned model for 1000 steps with no control perturbations
- The autocorrelation was compared to the actual Poincare map



Simplified EM: K-Means

- Rather than letting a training point x_j belong to several models, we force it to belong to only one
 - Choose the most likely model, $z_j = \operatorname{argmax}_{i=1,\dots,N} Q_j(Z = i; \bar{\mathbf{M}})$
 - Equivalent to using sparse $\bar{Q}_j^i = \begin{cases} 1 & i = z_j \\ 0 & i \neq z_j \end{cases}$
 - Speeds up regression and Gaussian mixture since fewer points are used to fit each model

1. Expectation Step: Find the most likely model at each training point

$$z_j \triangleq \operatorname{argmax}_{i=1,\dots,N} p_{Vi}(v_j^i; \bar{\mathbf{M}}) p_{X|Z=i}(x_j | Z = i; \bar{\mathbf{M}}) P(Z = i; \bar{\mathbf{M}})$$

- The collection of points assigned to each model will be used

$$S^i \triangleq \{j \in \{1, 2, \dots, m\} : z_j = i\}$$

2. Maximization Step: The usual maximization, but with sparse \bar{Q}_j^i

$$\bar{\mathbf{M}} \leftarrow \operatorname{argmax}_{\mathbf{M}} \sum_{i=1}^N \sum_{j \in S^i} \log[p_{Vi}(v_j^i; \mathbf{M}) p_{X|Z=i}(x_j | Z = i; \mathbf{M}) P(Z = i; \mathbf{M})]$$

K-Means Maximization Step

- Identical to EM maximization step, but with sparse weight matrices.

I. Regression: let $\mathbf{Y}^i = [\cdots \ y_j \ \cdots]_{j \in S^i}$ and $\mathbf{\Phi}^i = \begin{bmatrix} \theta_1(x_j, u_j) & & \\ \vdots & \ddots & \\ \theta_L(x_j, u_j) & & \end{bmatrix}_{j \in S^i}$

- We have simple least squares over a restricted set of points

$$A^i = (\mathbf{Y}^i \mathbf{\Phi}^{iT}) (\mathbf{\Phi}^i \mathbf{\Phi}^{iT})^+ \quad R^i = \frac{1}{|S^i|} (\mathbf{Y}^i - A^i \mathbf{\Phi}^i) (\mathbf{Y}^i - A^i \mathbf{\Phi}^i)^T$$

II. Gaussian Mixture:

- We fit the maximum likelihood Gaussian to each point cluster

$$\mu_x^i = \frac{1}{|S^i|} \sum_{j \in S^i} x_j \quad \Sigma^i = \frac{1}{|S^i|} \sum_{j \in S^i} (x_j - \mu_x^i)(x_j - \mu_x^i)^T$$

III. Multinoulli:

- Prior model probabilities are just the fraction of points in each cluster

$$\phi^i = \frac{1}{m} |S^i|$$

K-Means/EM Variants

Many variants are possible. For example:

1. Choose most likely model $z_j = \operatorname{argmax}_i Q_j(Z = i; \bar{\mathbf{M}})$ but keep the weights $\bar{Q}_j^i = \begin{cases} Q_j(Z = i; \bar{\mathbf{M}}) & i = z_j \\ 0 & i \neq z_j \end{cases} \rightarrow$ Weighted regression and Gaussian mixture.
2. Choose the top $N' < N$ models based on likelihood $Q_j(Z = i; \bar{\mathbf{M}})$ and zero-out the rest (call it “Generalized K-Means”)
3. Allocate training points by choosing among top $N' < N$ models based on minimum prediction error
 $z_j = \operatorname{argmin}_{i=1, \dots, N'} \|y_j - f^i(x_j, u_j)\|$
4. For each model (i) restrict to top $m' < m$ training points S^i and keep the weights on these points only (call it “Restricted EM”)

$$\bar{Q}_j^i = \begin{cases} Q_j(Z = i; \bar{\mathbf{M}}) & j \in S^i \\ 0 & j \notin S^i \end{cases}$$

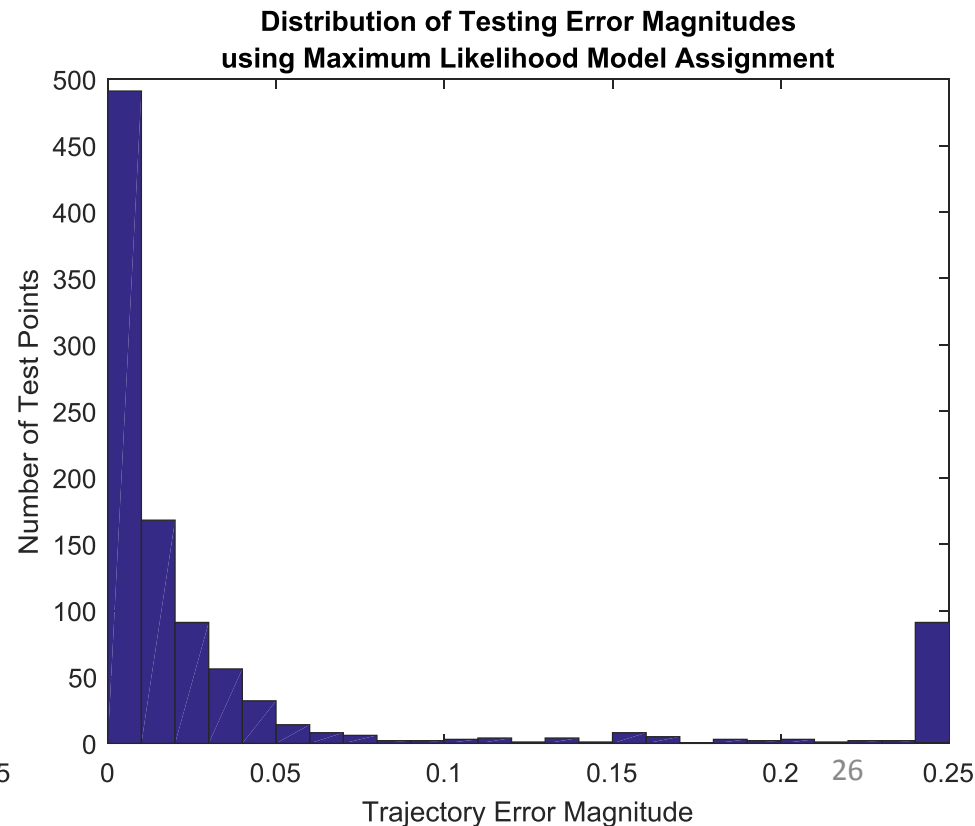
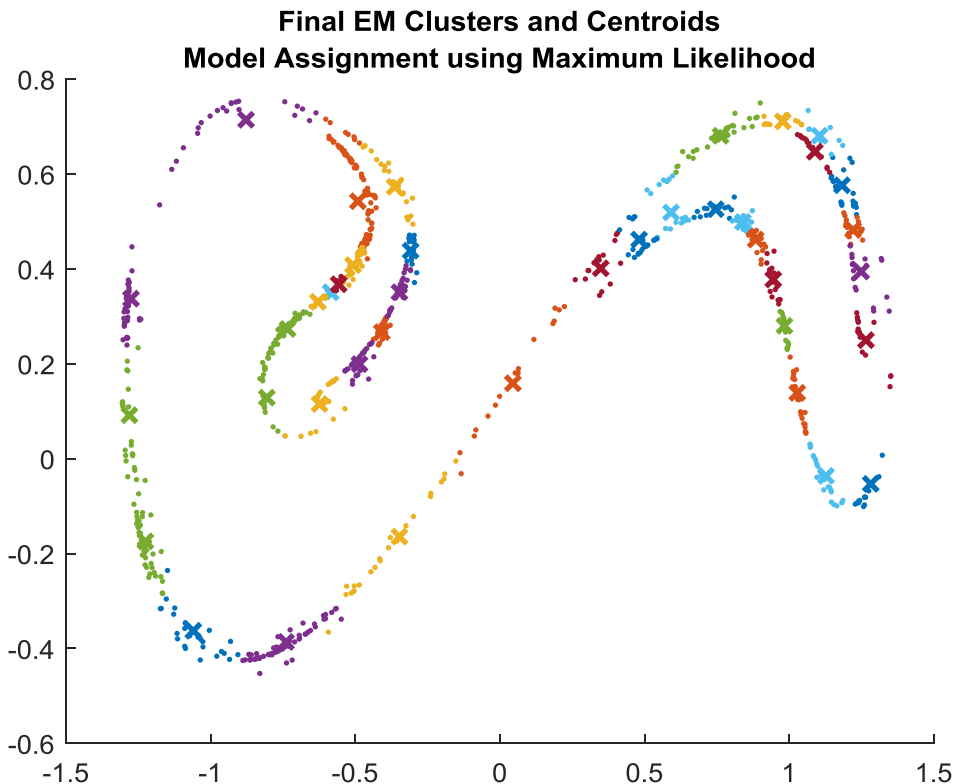
- These will enable efficient extensions to kernel regression.
- Intuitive since *our models only need to be accurate locally.*

Linear K-Means Gaussian Mixture Results

- $N = 40$ linear models were fit
- Each point is assigned to the single most likely model at each step
- Maximum likelihood model assignment is used to reconstruct the dynamics
 - Numerical experiments show that of the methods discussed, maximum likelihood produced the lowest errors.
- Minimum singular values of covariances Σ and R were needed to prevent degeneracy
- $m_{test} = 1000$ independently generated test points with random control perturbations were used for evaluation

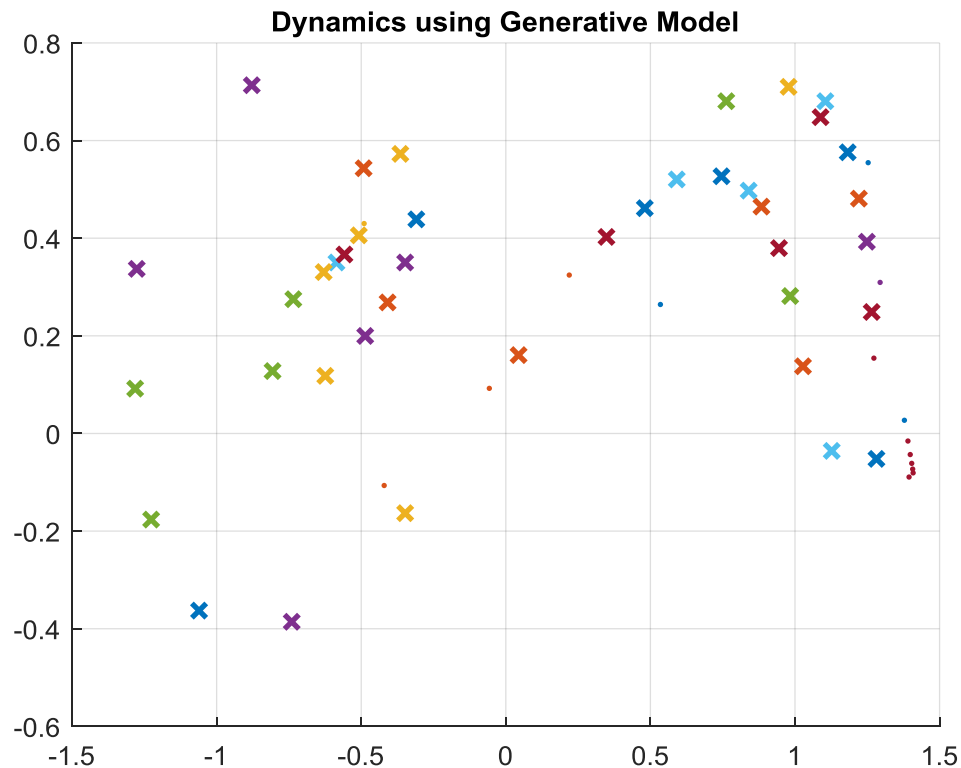
Linear K-Means Gaussian Mixture Results

- Mean Distance Error: $\bar{E}_D = 0.0842$
- Median Distance Error: 0.0103
- Decaying Error Norm over $M = 10$ steps with $\alpha = 2$: $\bar{E}_\alpha^M = 342.6198$



Use as a Generative Model

- The dynamics were simulated using the learned model for 1000 steps with no control perturbations
- The autocorrelation was compared to the actual Poincare map
- The approximation diverges and is therefore not useful as a generative model

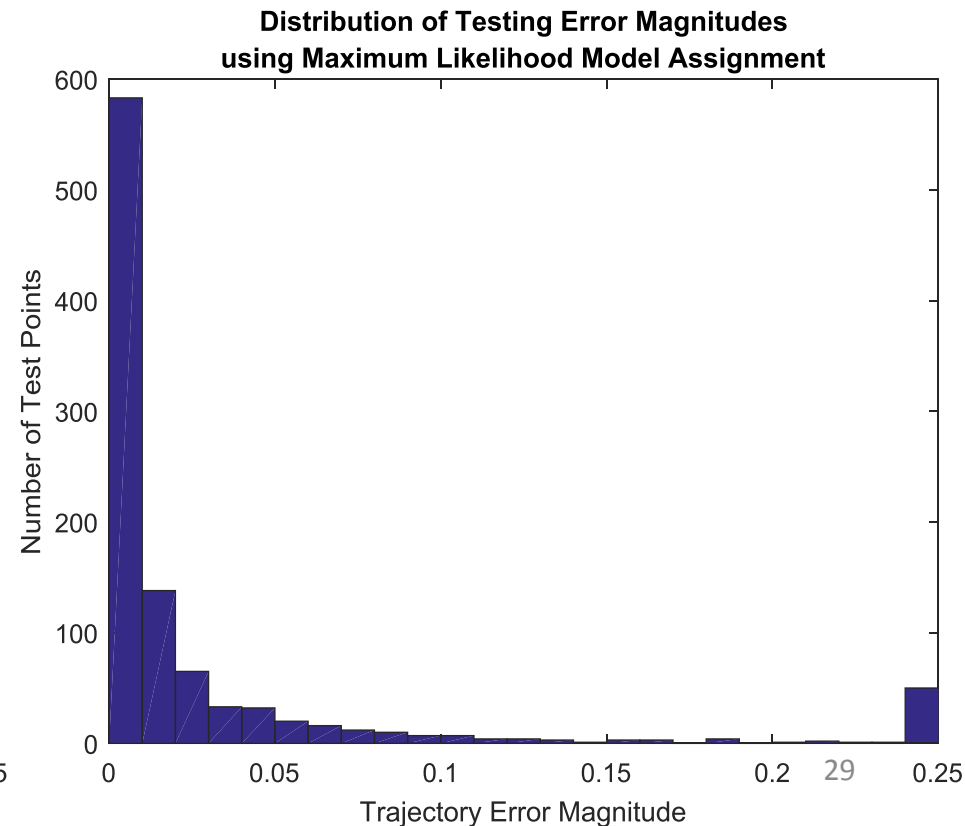
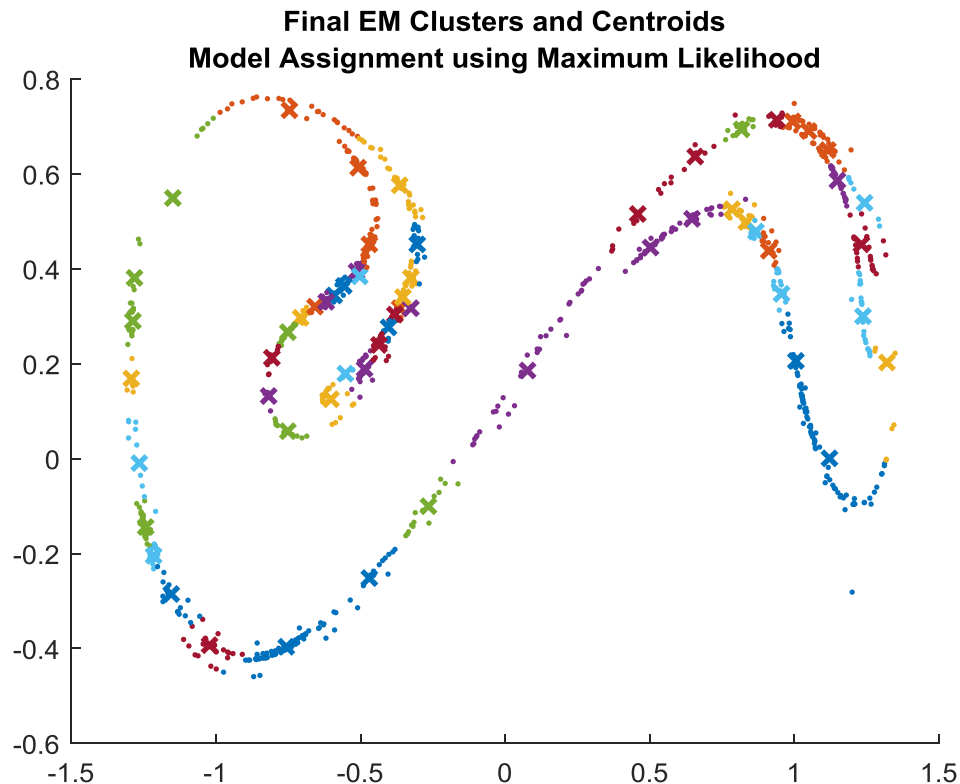


Linear Generalized K-Means Results

- $N = 60$ linear models were fit
- Each training point was assigned to $N' < N$ models based on maximum likelihood giving sparse \bar{Q}_j^i
- K-means was initially converged with $N' = 1$ then increased to $N' = 3$ and converged
 - This encourages the centroids to stay tight to the filamented distribution of points in the Poincare map
- Maximum likelihood model assignment is used to reconstruct the dynamics
- Minimum singular values of covariances Σ and R were needed to prevent degeneracy

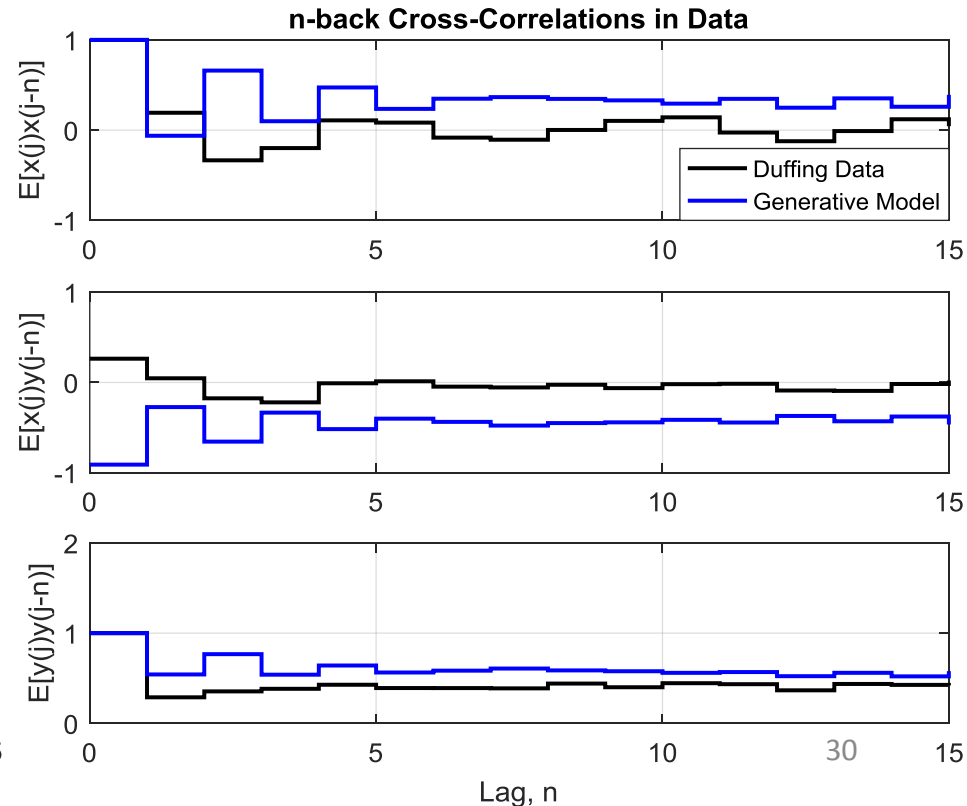
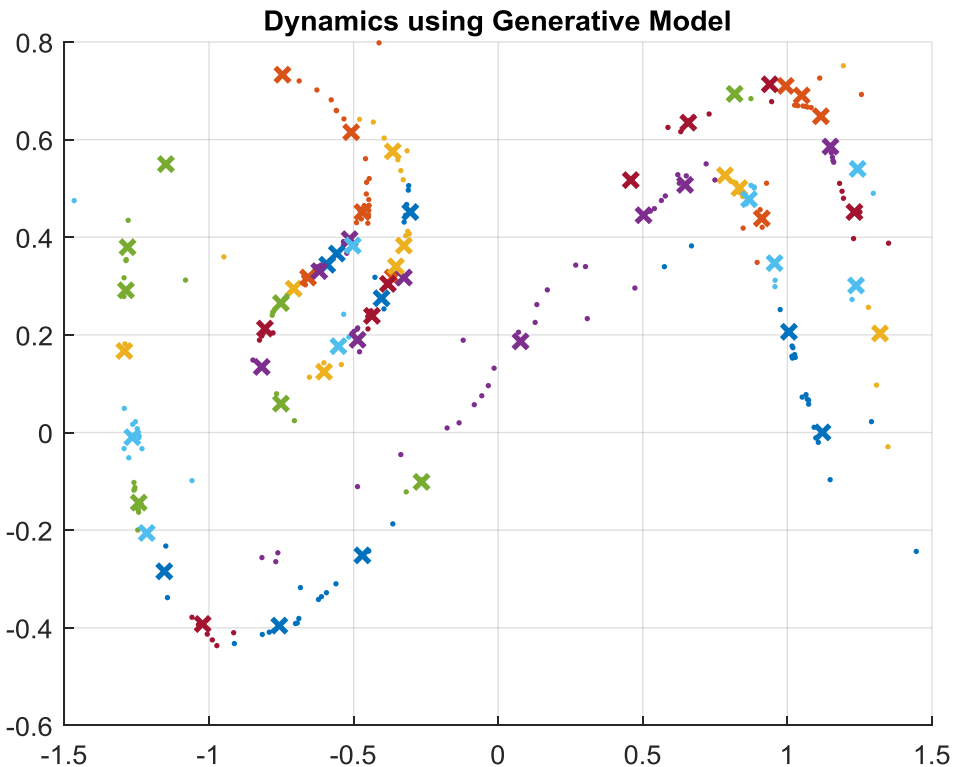
Linear Generalized K-Means Results

- Mean Distance Error: $\bar{E}_D = 0.0498$
- Median Distance Error: 0.0068
- Decaying Error Norm over $M = 10$ steps with $\alpha = 2$: $\bar{E}_\alpha^M = 0.5541$



Use as a Generative Model

- The dynamics were simulated using the learned model for 1000 steps with no control perturbations
- The autocorrelation was compared to the actual Poincare map



Highlights from Linear Cases

- Generalized K-means where points can be assigned to multiple models attains the lowest error with a large number of models
- The K-means and generalized K-means centroids fit the data more tightly than EM.
- K-means and generalized K-means both take less time to converge than EM even with the restricted set of points used to train each model.
- Linear models are very useful
 - Local DMD
- Linear models are still too inaccurate to replicate the dynamics over multiple steps
 - Nonlinear kernel models will now be developed

Extending to Nonlinear Local Dynamics

- As we have seen, the local dynamics can be highly nonlinear
 - The number of linear models that can be fit is limited by the number of training points
 - Cannot provide adequate resolution in areas with high nonlinearity unless they are accompanied by a high density of training point
- We would like to model nonlinear local variations efficiently without over-fitting.
- Must be able to include local weighting when fitting the nonlinear model
- These requirements suggest using *weighted kernel ridge regression*.

Weighted Kernel Ridge Regression

- Our solution to the regression problem posed probabilistically

$$\hat{y} = f^i(x, u) = (\mathbf{Y} \mathbf{W}^i \mathbf{\Phi}^T) (\mathbf{\Phi} \mathbf{W}^i \mathbf{\Phi}^T)^+ \mathbf{\Phi}(x, u)$$

Is the solution of an equivalent matrix optimization problem

$$\hat{y} = A^i \mathbf{\Phi}(x, u) \quad A^i = \operatorname{argmin}_{A \in \mathbb{R}^{n \times L}} \left\| (\mathbf{Y} - A \mathbf{\Phi}) \sqrt{\mathbf{W}^i} \right\|_F^2$$

- We introduce nonlinearity by considering a nonlinear feature map from the data $\xi = [x \quad u \quad 1]^T$ into a high-dimensional Hilbert space \mathcal{H}

$$\Phi: \mathbb{R}^{n+q+1} \rightarrow \mathcal{H}$$

$$\mathbf{\Phi}: \mathbb{R}^m \rightarrow \mathcal{H}$$

$$\xi \mapsto \Phi(\xi)$$

$$v \mapsto v_1 \Phi(\xi_1) + \cdots + v_m \Phi(\xi_m)$$

Weighted Kernel Ridge Regression

- Consider the l th component of the function we wish to approximate. We will find a (possibly infinite) combination $\psi_l \in \mathcal{H}$ of features in \mathcal{H} to represent the function

$$\hat{y}(l) = \langle \Phi(\xi), \psi_l \rangle_{\mathcal{H}} \approx f_l(x, u)$$

- For ease of notation, let

$$\underline{y}(l) \triangleq \begin{bmatrix} y_1(l) \\ \vdots \\ y_m(l) \end{bmatrix} \in \mathbb{R}^m \quad \langle \Phi, \psi_l \rangle_{\mathcal{H}} \triangleq \begin{bmatrix} \langle \Phi(\xi_1), \psi_l \rangle_{\mathcal{H}} \\ \vdots \\ \langle \Phi(\xi_m), \psi_l \rangle_{\mathcal{H}} \end{bmatrix} \in \mathbb{R}^m$$

- Then we pose the weighted linear regression problem in \mathcal{H} for the l th component of f

$$\psi_l^* = \operatorname{argmin}_{\psi_l \in \mathcal{H}} \left\| \sqrt{W^i} \left(\underline{y}(l) - \langle \Phi, \psi_l \rangle_{\mathcal{H}} \right) \right\|_2^2 + \lambda \langle \psi_l, \psi_l \rangle_{\mathcal{H}}$$

- We introduce the “Ridge” penalty $\lambda \langle \psi_l, \psi_l \rangle_{\mathcal{H}}$ to avoid overfitting the training data when there are an excess of nonlinear features

Learning Subspace Property (LSP)

- A.K.A. “Occam’s razor”
 - While there may be infinitely many $\psi_l \in \mathcal{H}$ satisfying the criteria, we look for the smallest $\psi_l \in \text{span}\{\Phi(\xi_j)\}_{j=1}^m = \mathcal{R}(\Phi)$ in the span of the training data
 - Analogous to right pseudoinverse of Φ^* in \mathcal{H}
- In the case of *ridge* regression, we can actually prove that $\psi_l^* \in \mathcal{R}(\Phi)$

Suppose that $\psi_l^* = \hat{\psi}_l + \tilde{\psi}_l$ where $\hat{\psi}_l \in \mathcal{R}(\Phi)$ and $\tilde{\psi}_l \in \mathcal{R}(\Phi)^\perp$

Then $\langle \Phi, \psi_l^* \rangle_{\mathcal{H}} = \langle \Phi, \hat{\psi}_l \rangle_{\mathcal{H}} + \langle \Phi, \tilde{\psi}_l \rangle_{\mathcal{H}} = \langle \Phi, \hat{\psi}_l \rangle_{\mathcal{H}}$

But $\langle \psi_l^*, \psi_l^* \rangle_{\mathcal{H}} = \langle \hat{\psi}_l, \hat{\psi}_l \rangle_{\mathcal{H}} + \langle \tilde{\psi}_l, \tilde{\psi}_l \rangle_{\mathcal{H}} \geq \langle \hat{\psi}_l, \hat{\psi}_l \rangle_{\mathcal{H}}$

$$\therefore \psi_l^* = \hat{\psi}_l \quad \text{and} \quad \exists a_l \in \mathbb{R}^m \quad \text{s.t.} \quad \psi_l^* = \Phi a_l$$

- The regression problem now becomes finite-dimensional

Weighted Kernel Ridge Regression

- Use LSP to express $\psi_l = \Phi a_l = a_{l_1} \Phi(\xi_1) + \dots + a_{l_m} \Phi(\xi_m)$

$$a_l^* = \operatorname{argmin}_{a_l \in \mathbb{R}^m} \left\| \sqrt{W^i} \left(\underline{y}(l) - \sum_{j=1}^m \langle \Phi, \Phi(\xi_j) \rangle_{\mathcal{H}} a_{l_j} \right) \right\|_2^2 + \lambda \sum_{i,j=1}^m a_{l_i} a_{l_j} \langle \Phi(\xi_i), \Phi(\xi_j) \rangle_{\mathcal{H}}$$

- Let $K \triangleq \langle \Phi, \Phi \rangle_{\mathcal{H}} \triangleq \left[\langle \Phi(\xi_i), \Phi(\xi_j) \rangle_{\mathcal{H}} \right]_{1 \leq i,j \leq m} \in \mathbb{R}^{m \times m}$ and do some simple manipulation

$$a_l^* = \operatorname{argmin}_{a_l \in \mathbb{R}^m} \left\| \sqrt{W^i} \left(\underline{y}(l) - K a_l \right) \right\|_2^2 + \lambda a_l^T K a_l$$

$$a_l^* = \operatorname{argmin}_{a_l \in \mathbb{R}^m} \left\| \begin{bmatrix} \sqrt{W^i} \underline{y}(l) \\ 0 \end{bmatrix} - \begin{bmatrix} \sqrt{W^i} K \\ \sqrt{\lambda K} \end{bmatrix} a_l \right\|_2^2$$

- The above is now a standard m -dimensional least squares problem whose solution is

$$a_l^* = (W^i K + \lambda I_m)^{-1} W^i \underline{y}(l)$$

$$\therefore \hat{y}(l) = \langle \Phi(\xi), \psi_l^* \rangle_{\mathcal{H}} = \langle \Phi(\xi), \Phi \rangle_{\mathcal{H}} (W^i K + \lambda I_m)^{-1} W^i \underline{y}(l)$$

Weighted Kernel Ridge Regression

- Notice that the entire solution to the regression problem is expressed in terms of inner products of elements in \mathcal{H}
 - Introduce kernel function $k(\xi, \eta) \triangleq \langle \Phi(\xi), \Phi(\eta) \rangle_{\mathcal{H}} \in \mathbb{R}$ and $\mathbf{k}(\xi) \triangleq \langle \Phi, \Phi(\xi) \rangle_{\mathcal{H}} \in \mathbb{R}^m$

- Then the full estimator for the function is given by

$$\hat{y} = f^i(x, u) = \mathbf{Y} \mathbf{W}^i (\mathbf{K} \mathbf{W}^i + \lambda \mathbf{I}_m)^{-1} \mathbf{k}(\xi) \quad \mathbf{K} = [k(\xi_i, \xi_j)]_{1 \leq i, j \leq m} \quad \xi = \begin{bmatrix} x \\ u \\ 1 \end{bmatrix}$$

- This approximation is easily differentiated if we know the derivative of the kernel with respect to one of its slots $D_{\xi} k(\xi_0, \xi)$ allowing us to find $D_{\xi} \mathbf{k}(\xi_0) \in \mathbb{R}^{m \times m}$

$$D_{\xi} f^i(\xi_0) \Delta \xi = \mathbf{Y} \mathbf{W}^i (\mathbf{K} \mathbf{W}^i + \lambda \mathbf{I}_m)^{-1} D_{\xi} \mathbf{k}(\xi_0) \Delta \xi$$

Kernel Regression Implementation

- It is prohibitively expensive to invert the $(KW^i + \lambda I_m) \in \mathbb{R}^{m \times m}$ each time the models are updated
- We therefore choose only the largest $m' < m$ weights $\{\bar{Q}_{j_1}^i, \bar{Q}_{j_2}^i, \dots, \bar{Q}_{j_{m'}}^i\}$ on training points $S^i = \{j_1, j_2, \dots, j_{m'}\}$ to be nonzero. \rightarrow “Restricted EM”
 - With the intention of doing *local* modeling, this is a sensible thing to do to speed up the algorithm
- Forming the sub-matrices

$$\tilde{K} = [K_{i_r j_s}]_{1 \leq r, s \leq m'} \quad \tilde{W} = \text{diag}\{\bar{Q}_{j_r}^i\}_{1 \leq r \leq m'} \quad \tilde{Y} = [y_{j_1} \cdots y_{j_{m'}}]$$

$$\tilde{\mathbf{h}}(\xi) = \left[\langle \Phi(\xi_{j_1}), \psi_l \rangle_{\mathcal{H}} \quad \cdots \quad \langle \Phi(\xi_{j_{m'}}), \psi_l \rangle_{\mathcal{H}} \right]^T$$

Using the Schur complement, we find that

$$\hat{y} = f^i(x, u) = \tilde{Y} \tilde{W}^i (\tilde{K} \tilde{W}^i + \lambda I_{m'})^{-1} \tilde{\mathbf{h}}(\xi)$$

Application to Duffing Equation with Control (Restricted EM Algorithm)

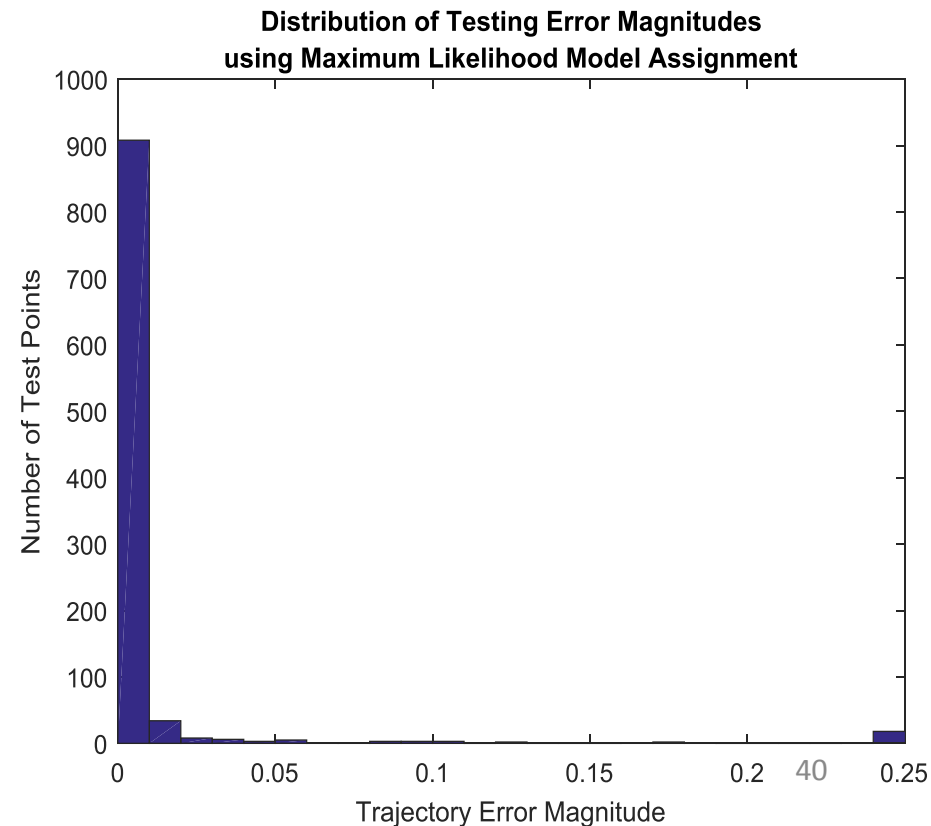
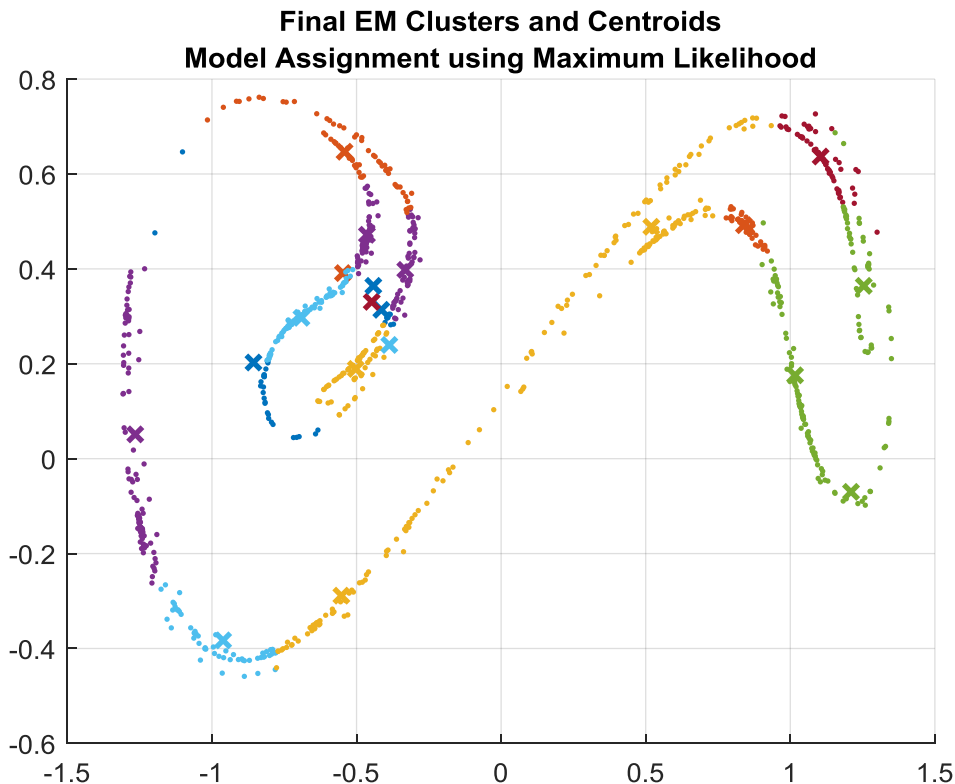
- $N = 20$ models were fit using the RBF kernel with $\sigma = 0.1$ and small $\lambda = 5 * 10^{-4}$

$$k(x, z) = \exp\left(-\frac{1}{2\sigma^2} \|x - z\|_2^2\right)$$

- Each model is fit to the $m' = 500$ training points with the largest weights \bar{Q}_j^i
- Maximum likelihood model assignment was used to reconstruct the dynamics
 - Numerical experiments show that of the methods discussed, maximum likelihood produced the lowest errors.
- Minimum singular values of covariances Σ and R were needed to prevent degeneracy
- $m_{test} = 1000$ independently generated test points with random control perturbations were used for evaluation

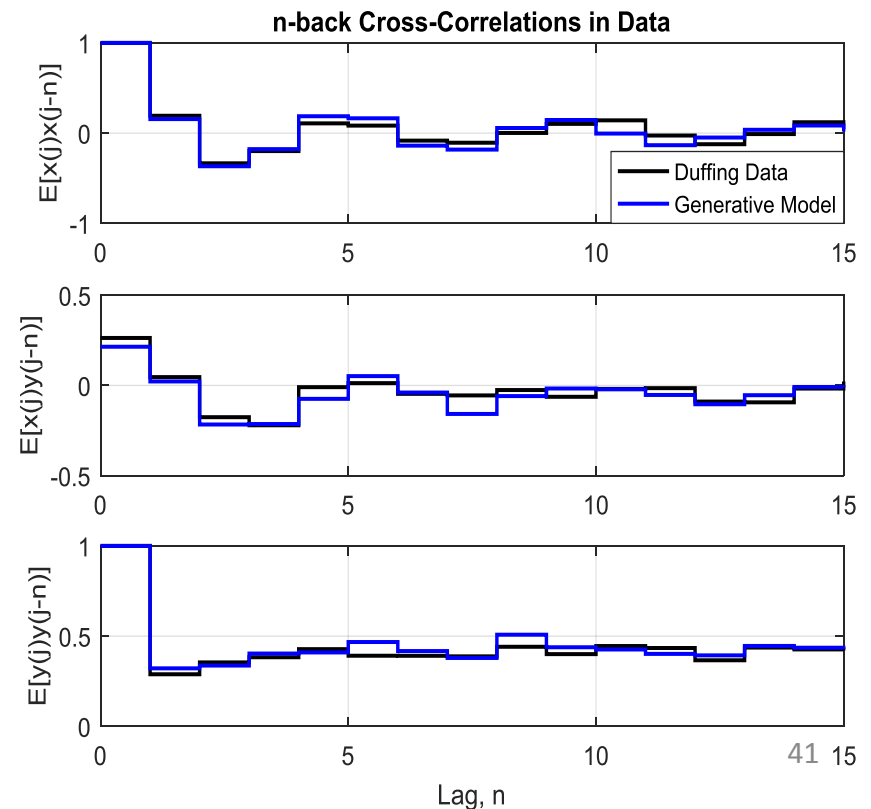
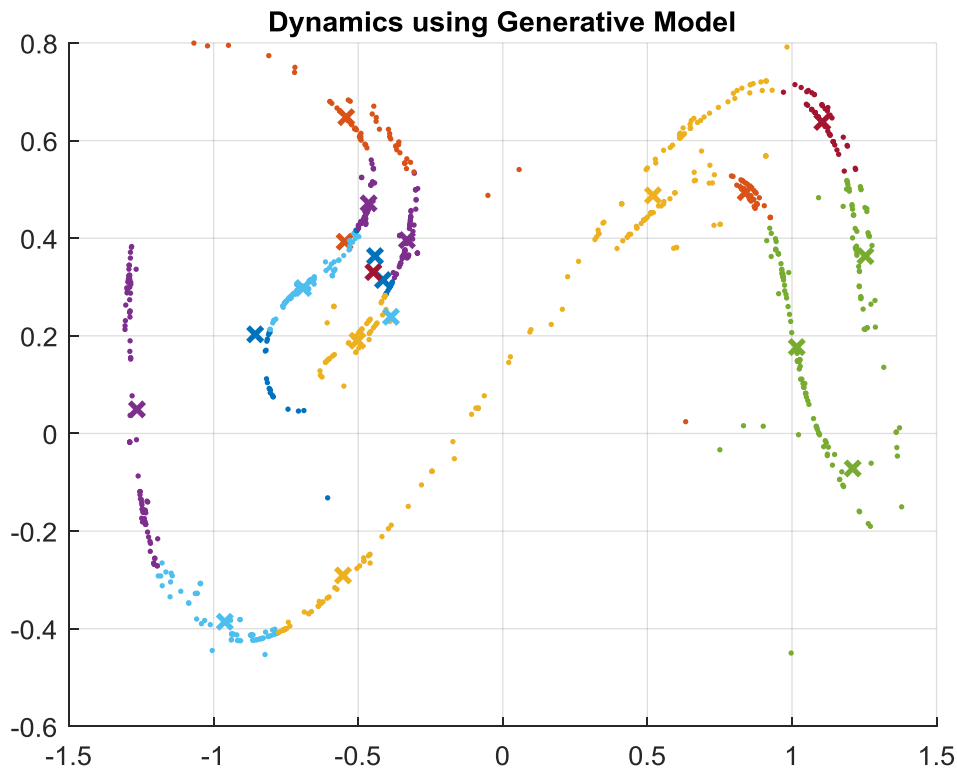
Nonlinear Restricted EM Results

- Mean Distance Error: $\bar{E}_D = 0.0209$
- Median Distance Error: 0.0014
- Decaying Error Norm over $M = 10$ steps with $\alpha = 2$: $\bar{E}_\alpha^M = 0.1196$



Use as a Generative Model

- The dynamics were simulated using the learned model for 1000 steps with no control perturbations
- The autocorrelation was compared to the actual Poincare map



Application to Duffing Equation with Control (K-Means Algorithm)

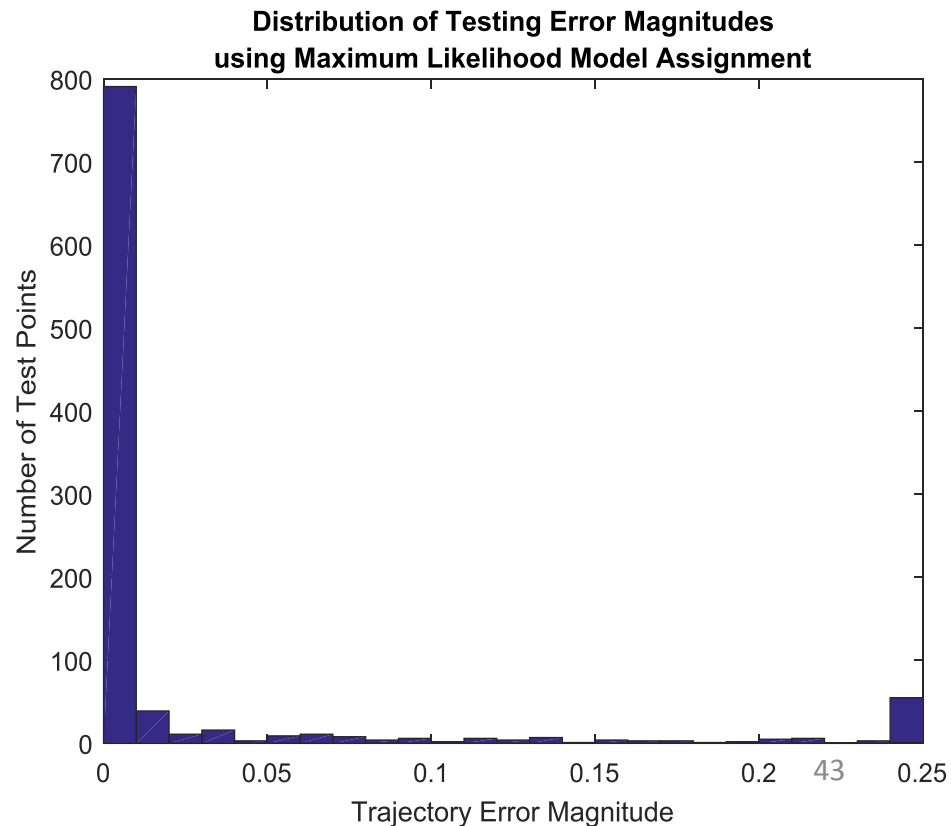
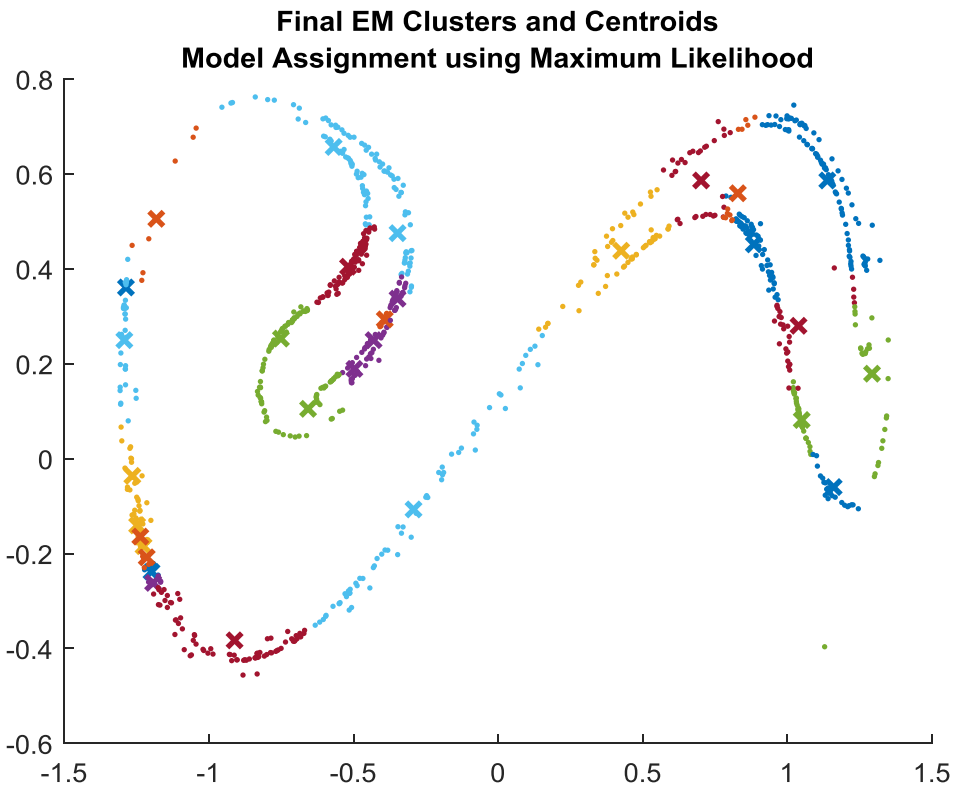
- $N = 30$ models were fit using the RBF kernel with $\sigma = 0.1$ and small $\lambda = 5 * 10^{-4}$

$$k(x, z) = \exp\left(-\frac{1}{2\sigma^2} \|x - z\|_2^2\right)$$

- Each training point was assigned to a single model based on maximum likelihood giving sparse \bar{Q}_j^i
 - Convergence was much faster even though more models were fit
- Maximum likelihood model assignment was used to reconstruct the dynamics
- Minimum singular values of covariances Σ and R were needed to prevent degeneracy
- $m_{test} = 1000$ independently generated test points with random control perturbations were used for evaluation

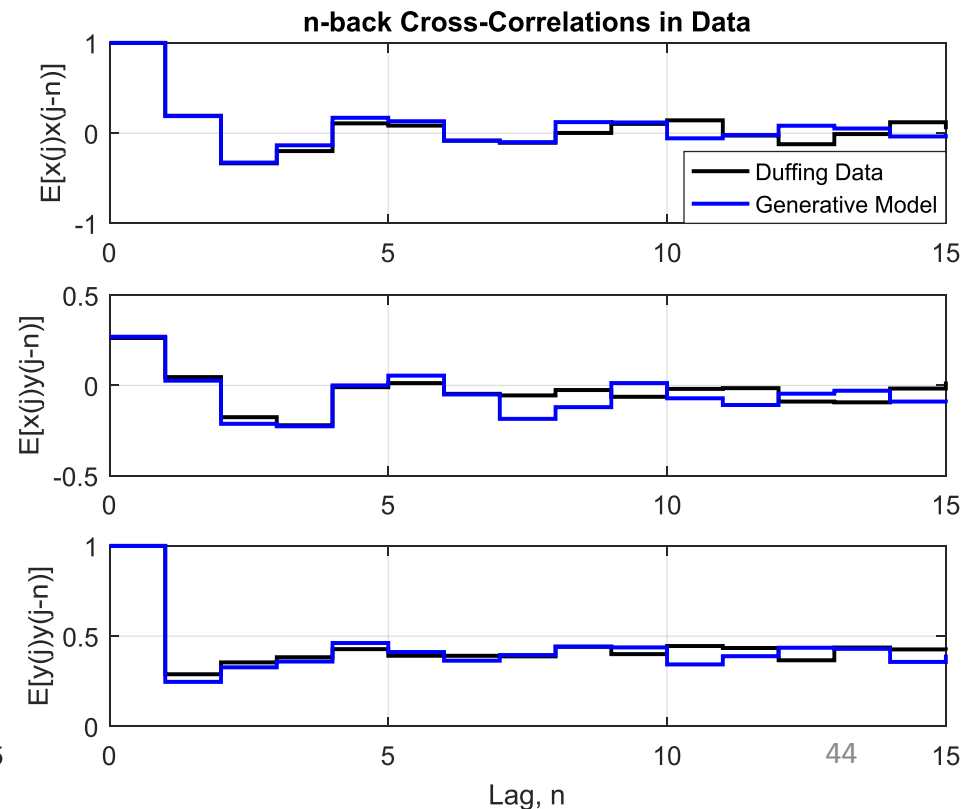
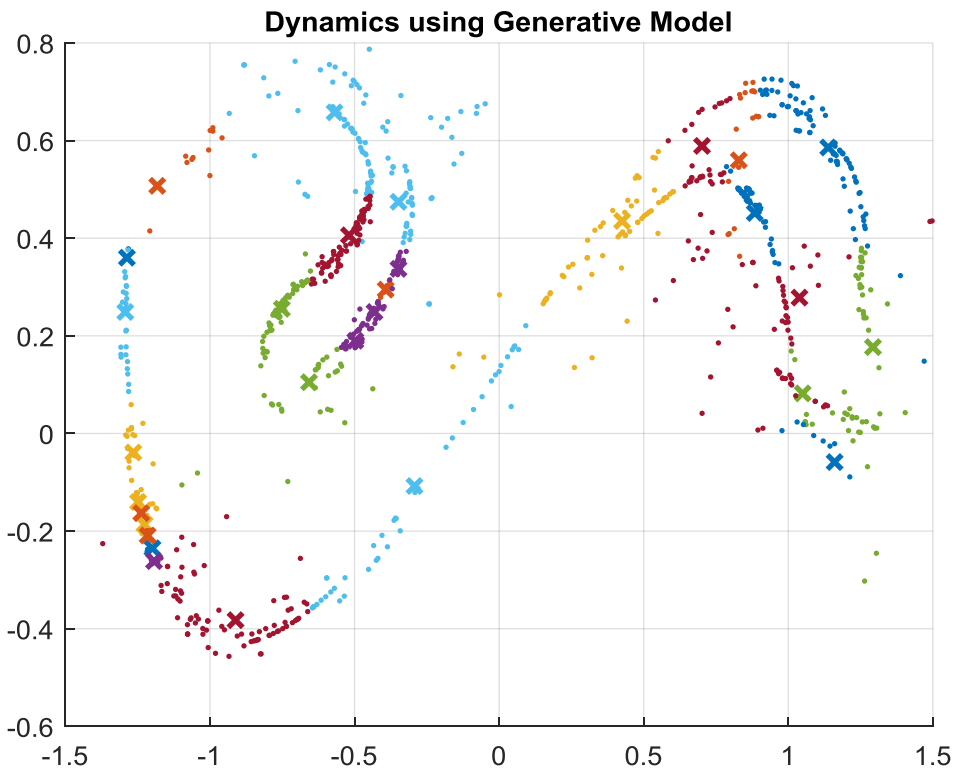
Nonlinear K-Means Results

- Mean Distance Error: $\bar{E}_D = 0.0500$
- Median Distance Error: 0.0021
- Decaying Error Norm over $M = 10$ steps with $\alpha = 2$: $\bar{E}_\alpha^M = 0.2765$



Use as a Generative Model

- The dynamics were simulated using the learned model for 1000 steps with no control perturbations
- The autocorrelation was compared to the actual Poincare map



Application to Duffing Equation with Control (Generalized K-Means Algorithm)

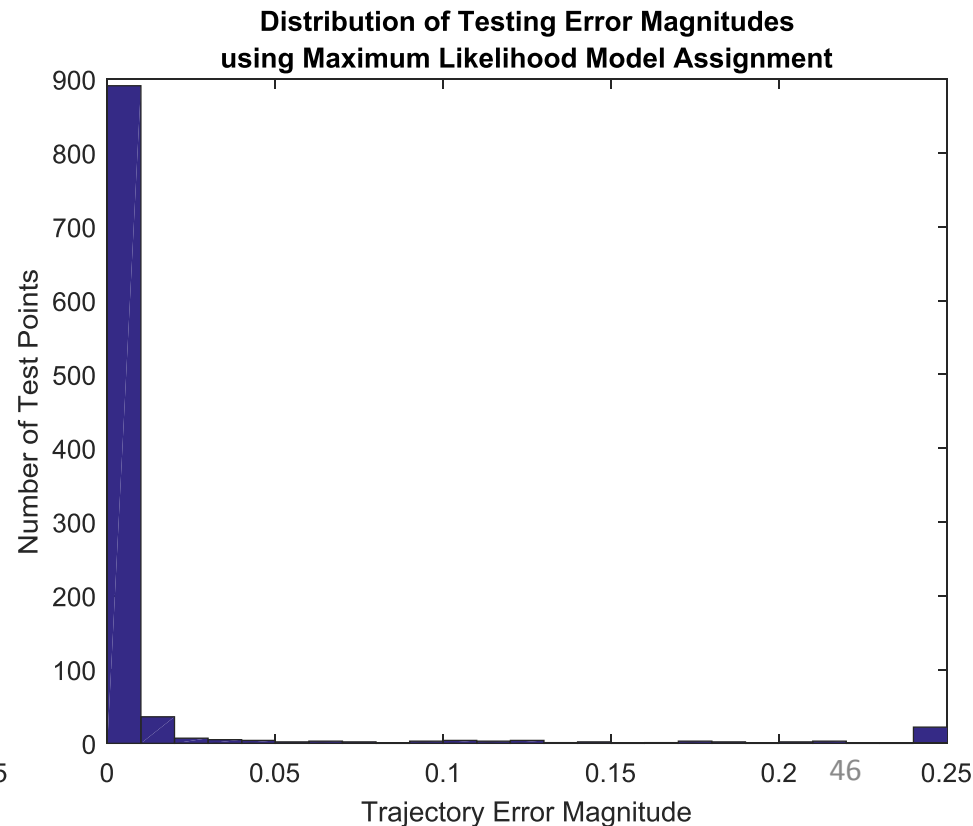
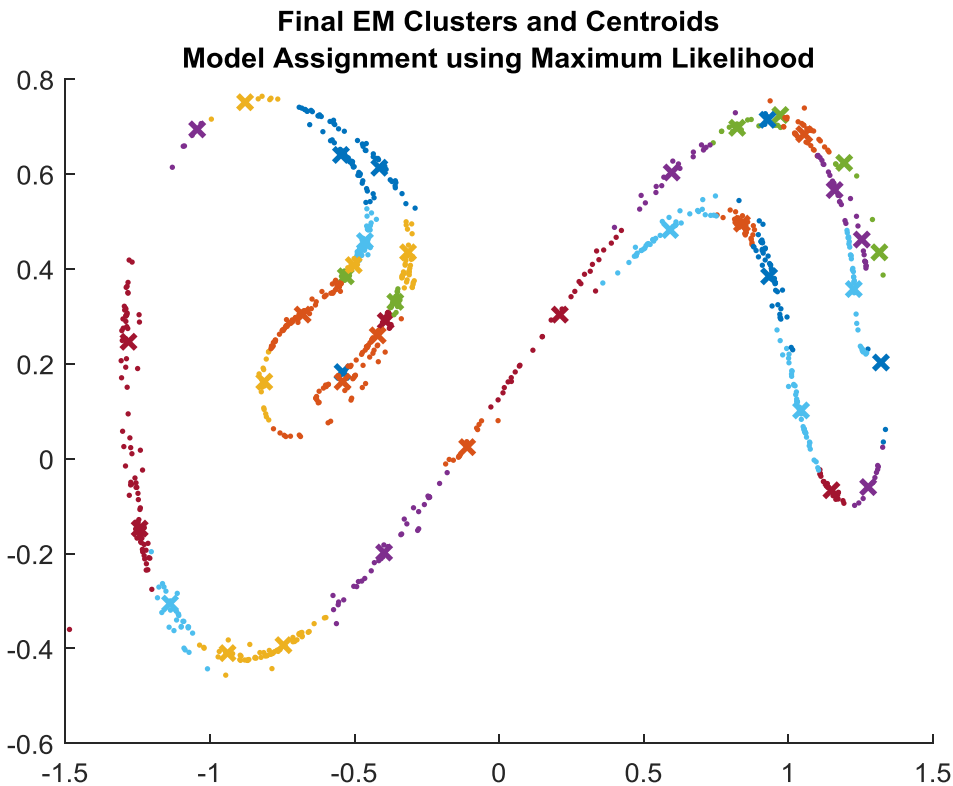
- $N = 40$ models were fit using the RBF kernel with $\sigma = 0.1$ and small $\lambda = 5 * 10^{-4}$

$$k(x, z) = \exp\left(-\frac{1}{2\sigma^2} \|x - z\|_2^2\right)$$

- Each training point was assigned to $N' < N$ models based on maximum likelihood giving sparse \bar{Q}_j^i
- K-means was initially converged with $N' = 1$ then increased to $N' = 4$ and converged
 - This encourages the centroids to stay tight to the filamented distribution of points in the Poincare map
- Maximum likelihood model assignment was used to reconstruct the dynamics
- Minimum singular values of covariances Σ and R were needed to prevent degeneracy
- $m_{test} = 1000$ independently generated test points with random control perturbations were used for evaluation

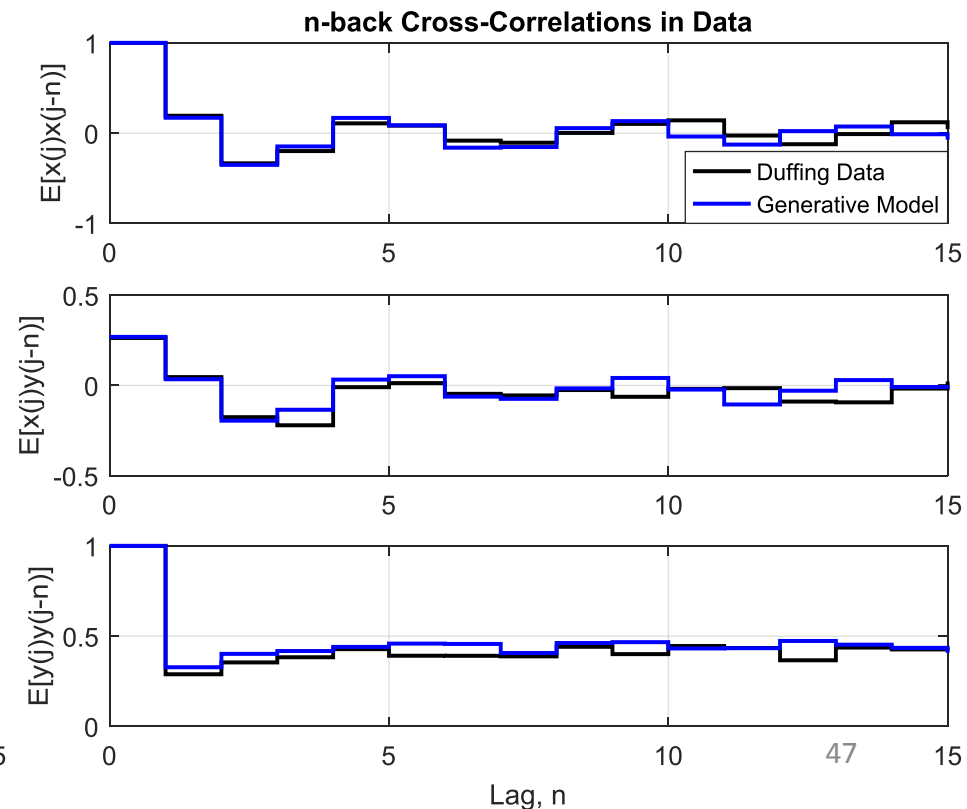
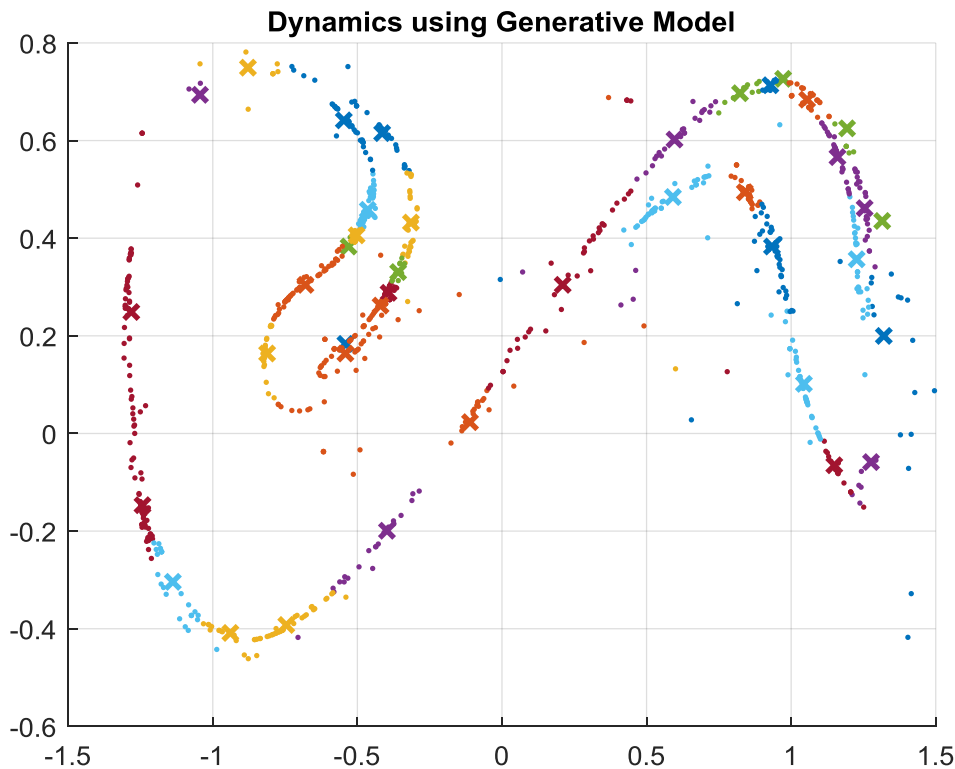
Nonlinear Generalized K-Means Results

- Mean Distance Error: $\bar{E}_D = 0.0208$
- Median Distance Error: 0.0017
- Decaying Error Norm over $M = 10$ steps with $\alpha = 2$: $\bar{E}_\alpha^M = 0.1434$



Use as a Generative Model

- The dynamics were simulated using the learned model for 1000 steps with no control perturbations
- The autocorrelation was compared to the actual Poincare map



Highlights from Nonlinear Cases

- The Restricted EM algorithm has the lowest error with the fewest number of models.
 - The algorithm takes a long time to converge
 - Centroids are not tight to filamented point distribution
- K-means converges quickly when points are assigned to only a single model.
 - Higher error
 - More models are required
 - Centroids are tight to filamented point distribution
- Generalized K-means where points are assigned to multiple models improves the error over K-means.
 - Many models are required
 - Low error
 - Centroids are tight to filamented point distribution

Non-Gaussian Local Distributions

- We have shown that fitting nonlinear kernel models allows us to better capture the dynamics in our model
- However, The Gaussian densities themselves do not reflect the local curvature and filamentation of the Chaotic Poincare map
 - Therefore, many models are needed
- It would be helpful to allow the local distributions $p_{X|Z}(x|Z)$ to become curved and stretched
- This is a work in progress, but it can be accomplished by fitting Gaussian densities in a kernel Hilbert space

$$p_{X|Z}(x|Z) = C \exp \left[-\frac{1}{2} \langle \tilde{\Sigma}(\Phi(x) - \mu_{\Phi}), (\Phi(x) - \mu_{\Phi}) \rangle_{\mathcal{H}} \right]$$
$$\mu_{\Phi} = \frac{1}{|S^i|} \sum_{j \in S^i} \Phi(x_j)$$

Gaussians in Feature Space

- Roughly the idea is to find the covariance of the data $\{x_j\}_{j=1}^m$ in feature space \mathcal{H} where

$$\Phi: \mathbb{R}^m \rightarrow \mathcal{H}, \quad v \mapsto v_1 \Phi(x_1) + \cdots + v_m \Phi(x_m)$$

$$\Sigma = \frac{1}{m} \left(\Phi - \frac{1}{m} \Phi \mathbf{1}_m \mathbf{1}_m^T \right) \left(\Phi - \frac{1}{m} \Phi \mathbf{1}_m \mathbf{1}_m^T \right)^* : \mathcal{H} \rightarrow \mathcal{H}$$

$$\Sigma = \frac{1}{m} \Phi \left(I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T \right) \Phi^*$$

- Then guess the form of the inverse operator $\tilde{\Sigma}$ restricted to the learning subspace $\mathcal{R}(\Phi)$
 - We know that its domain and range are $\mathcal{R}(\Phi)$
 - We know that it is self-adjoint

$$\tilde{\Sigma} = \Phi B \Phi^* \quad \text{for some} \quad B = B^T \in \mathbb{R}^{m \times m}$$

Gaussians in Feature Space

- In order for $\tilde{\Sigma} = \Phi B \Phi^*$ to be the inverse of $\Sigma = \frac{1}{m} \Phi \left(I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T \right) \Phi^*$ over the learning subspace, we must have

$$K = \Phi^* \Phi = \Phi^* \tilde{\Sigma} \Sigma \Phi = \frac{1}{m} \underbrace{\Phi^* \Phi}_K B \underbrace{\Phi^* \Phi}_K \left(I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T \right) \underbrace{\Phi^* \Phi}_K$$

And

$$K = \Phi^* \Phi = \Phi^* \Sigma \tilde{\Sigma} \Phi = \frac{1}{m} \underbrace{\Phi^* \Phi}_K \left(I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T \right) \underbrace{\Phi^* \Phi}_K B \underbrace{\Phi^* \Phi}_K$$

- We must have

$$B = m \left[K \left(I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T \right) K \right]^+$$

- Therefore, the non-normalized density is given by

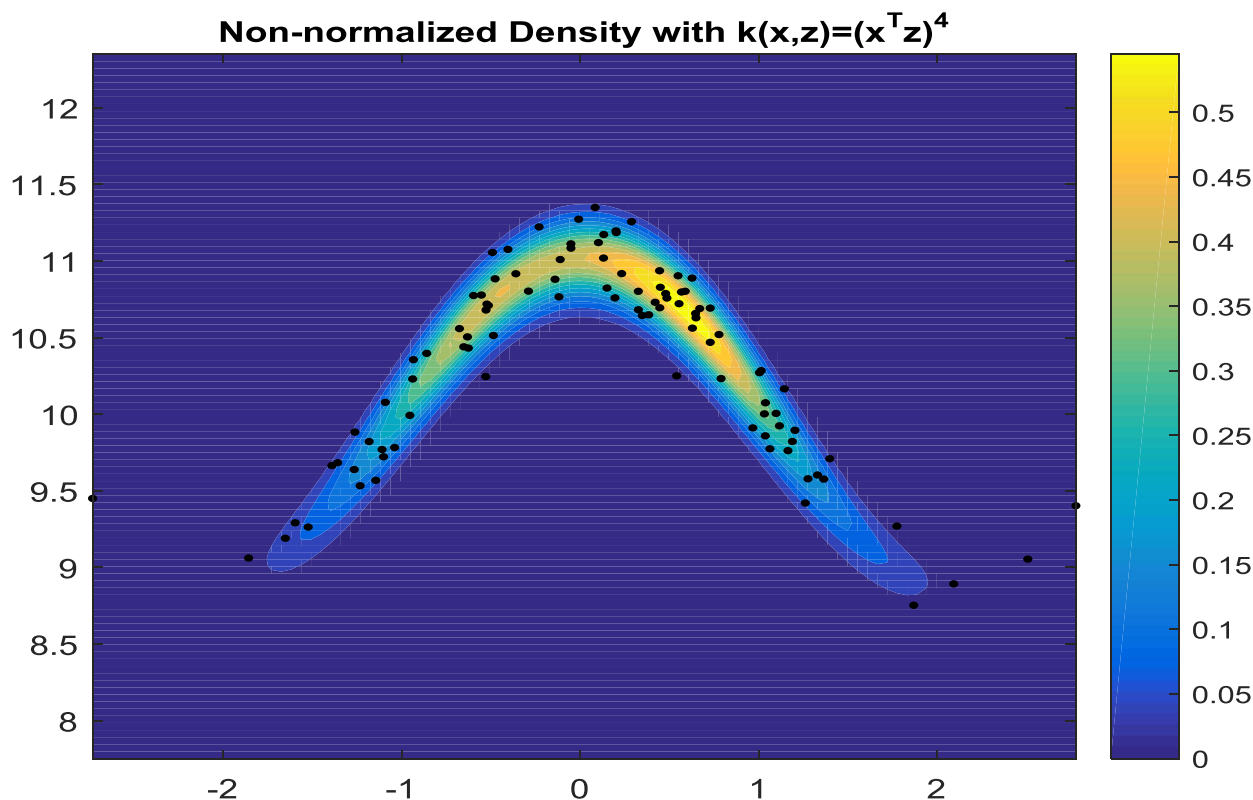
$$p_{X|Z}(x|Z) = C \exp \left[-\frac{1}{2} \left(\Phi(x) - \frac{1}{m} \Phi \mathbf{1}_m \right)^* \Phi B \Phi^* \left(\Phi(x) - \frac{1}{m} \Phi \mathbf{1}_m \right) \right]$$

$$p_{X|Z}(x|Z) = C \exp \left[-\frac{1}{2} \left(\boldsymbol{\kappa}(x) - \frac{1}{m} K \mathbf{1}_m \right)^* B \left(\boldsymbol{\kappa}(x) - \frac{1}{m} K \mathbf{1}_m \right) \right]$$

- I don't know of a good direct way to find C , but we do have a very convenient Monte-Carlo sample $\{x_1, \dots, x_m\}$ to approximate it!

Gaussians in Feature Space

- (Macready, W. G., “Density Estimation with Mercer Kernels,” NASA Ames Research Center) presents an EM algorithm for fitting mixtures of Gaussians in nonlinear feature space.
- The following figure was made using the method described above with $m = 100$ points and a 4th order polynomial kernel. The RBF kernel tends to over-fit.



Many Avenues for Future Work

- Use linear models to perform local DMD in identified regions of interest
 - Spectral analysis is also possible using the kernel model's Jacobian $D_{\xi} f^i(\xi_0)$
- Identifying global topological structure
 - Maximum likelihood model classification partitions the space efficiently → Look at Markov transition probabilities between models
 - Patch together Gaussian “pancakes” to study the underlying data manifold → non-parametric manifold learning (look at Lukaszuk-Karmowski metric)
 - Help identify periodic orbits, tori, and attractors

Many Avenues for Future Work

- Use linear models and conditional probabilities to do state estimation
 - Account for noisy observation process
 - Develop Kalman filter
- Possibility for online formulation
 - New data can be used to update models using matrix inversion lemma
 - K-means clustering can be performed “greedily”
 - Application to streaming DMD
- Given that both parts (Regression and Gaussian Mixture fitting) of the maximization step can be kernelized suggests that *the entire algorithm might be performed in high-dimensional intrinsic space and kernelized!*

Many Avenues for Future Work

- Investigate how the algorithm behaves for high-dimensional systems
 - Dealing with degenerate Gaussians and dynamics
 - Learning data manifold and locally projecting onto it
 - Investigate dimensionality reduction (kernel PCA or autoencoder) prior to model fitting
- Continuous reconstruction of the model estimates also require further investigation
 - Need some kind of probability-weighted smooth interpolation
- Implementation note: Since each model in the Maximization step is fit independently, the algorithm may be executed in parallel.