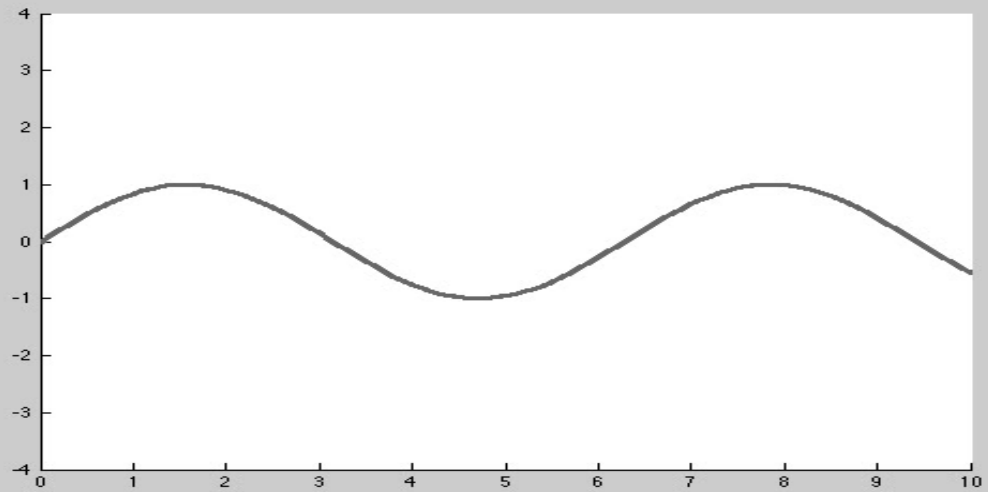


CSE-571:

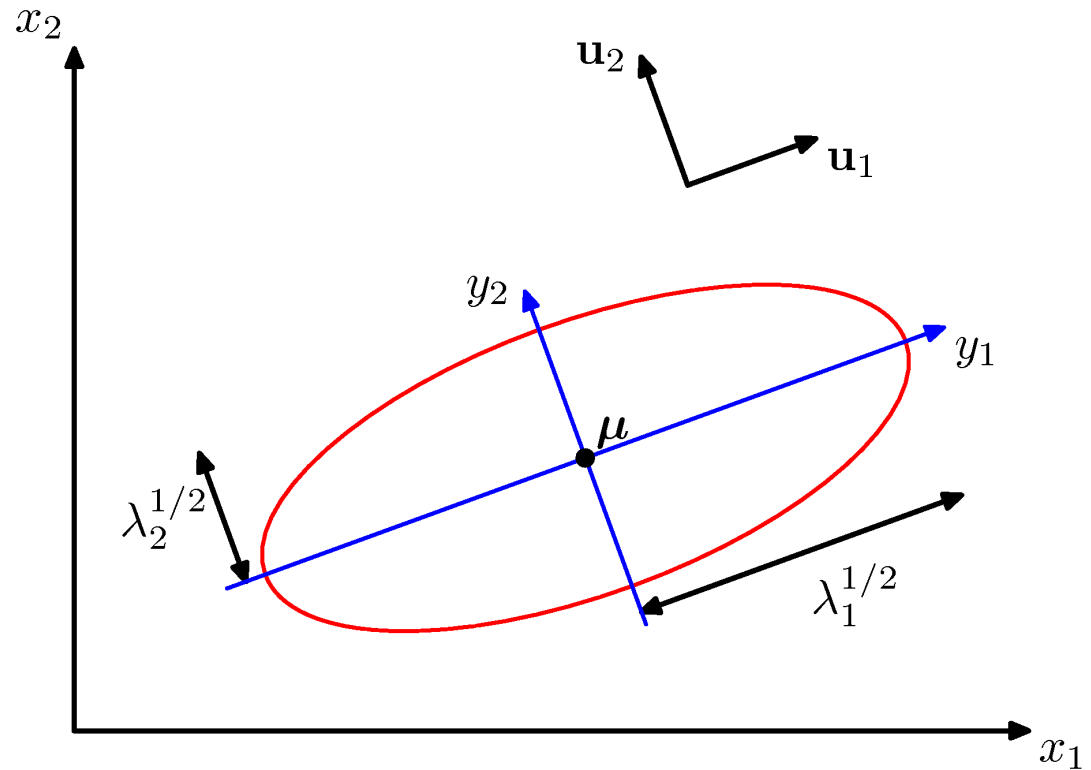
Gaussian Processes for Bayesian Filtering



High-level Idea of GPs

- Non-parametric regression model
- Distribution over functions
- Fully specified by training data and kernel function
- Output variables are jointly Gaussian
- Covariance given by distance of inputs in kernel space

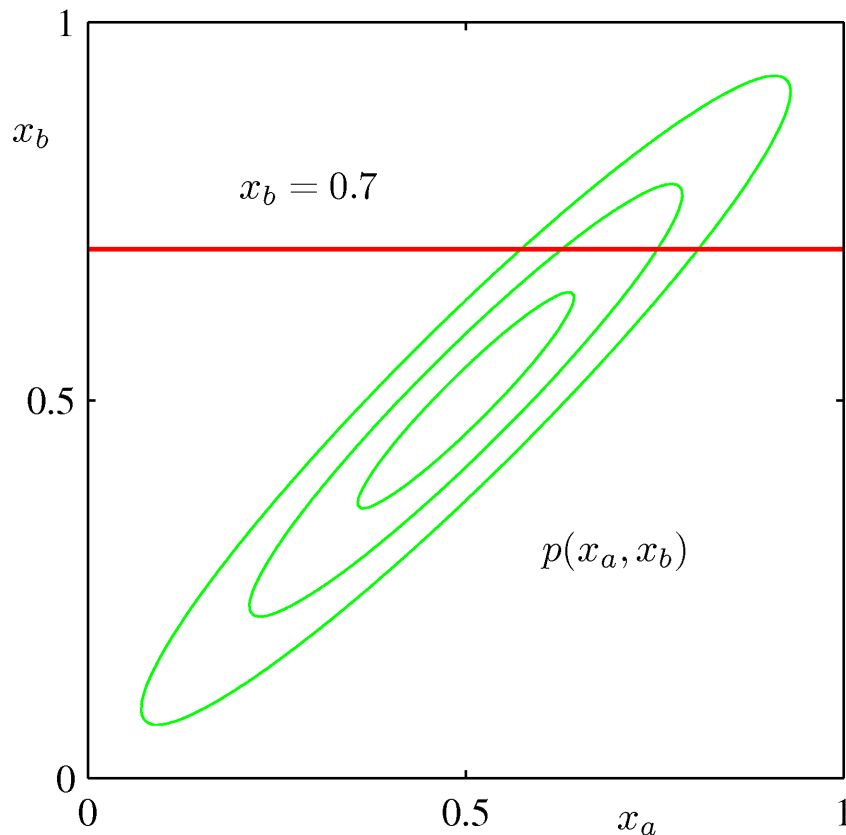
Gaussians



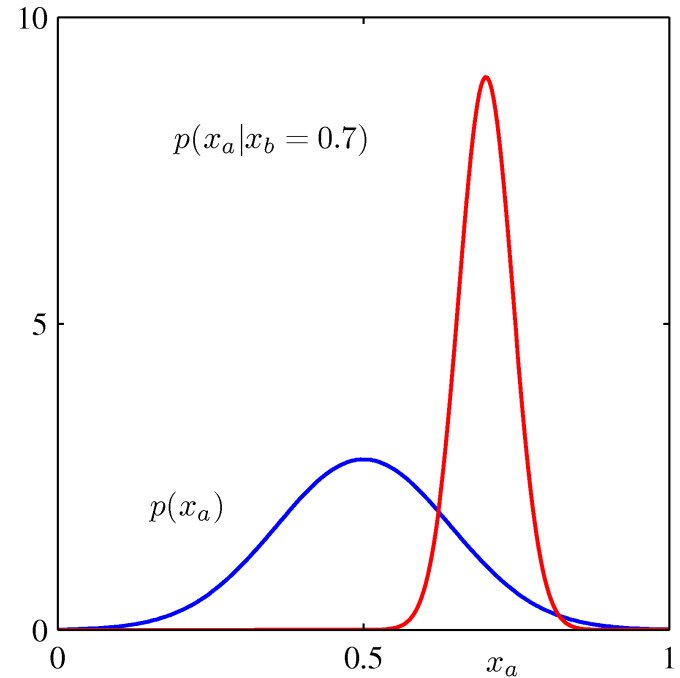
$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

Marginalization / Conditioning



$$p(x_a) = \mathbf{N}(\mu_a, \Sigma_{aa})$$



$$p(x_a | x_b) = \mathbf{N}(\mu_{a|b}, \Sigma_{a|b})$$

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

GP Setting

- Outputs are noisy function of inputs:

$$y_i = f(\mathbf{x}_i) + \varepsilon$$

- GP prior: Outputs jointly zero-mean Gaussian:

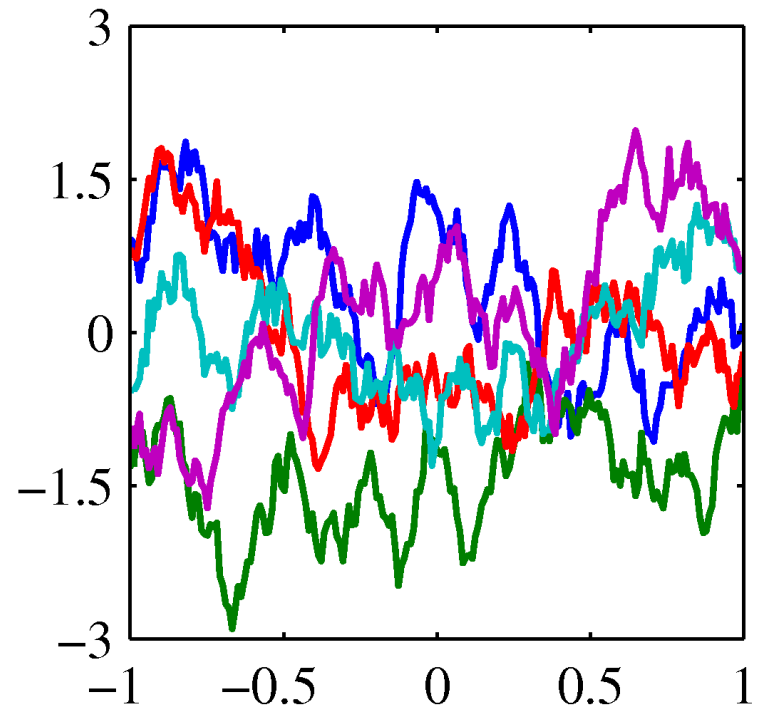
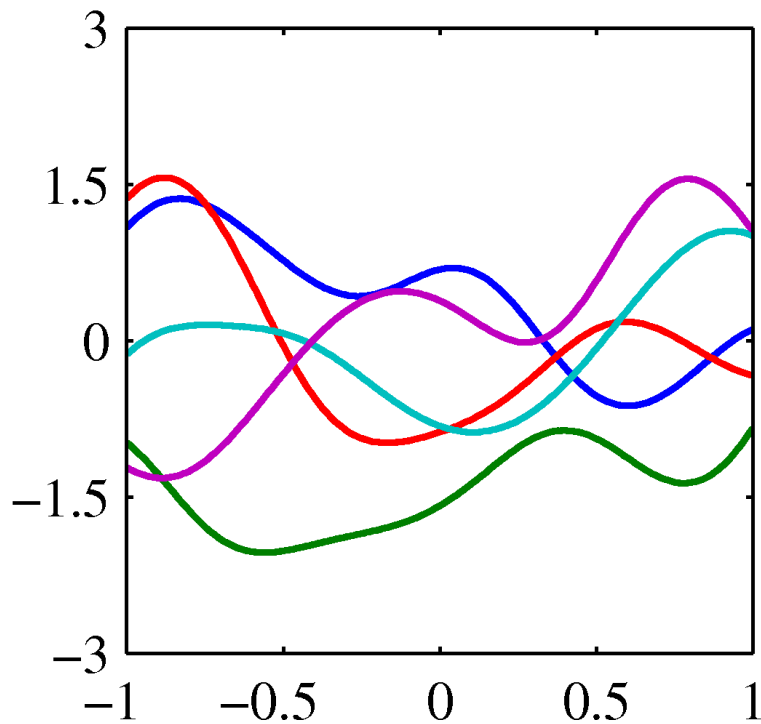
$$p(\mathbf{y} | \mathbf{X}) = \mathbf{N}(\mathbf{0}, \mathbf{K} + \sigma_n^2 \mathbf{I})$$

- Covariance given by kernel matrix over inputs:

$$\mathbf{K} = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & & \vdots \\ \vdots & k(x_i, x_i) & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix}$$

$$k(x, x') = \sigma_f^2 e^{-\frac{1}{2}(x-x')^T W (x-x')}$$

Functions Sampled from Prior



GP Prediction

- Training data:

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} = (\mathbf{X}, \mathbf{y})$$

- Prediction given training samples:

$$p(y^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) = N(\mu, \sigma^2)$$

$$\mu = \mathbf{k}^{*T} (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbf{k}^*[i] = k(\mathbf{x}^*, \mathbf{x}_i)$$

$$\sigma^2 = k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}^{*T} (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}^*$$

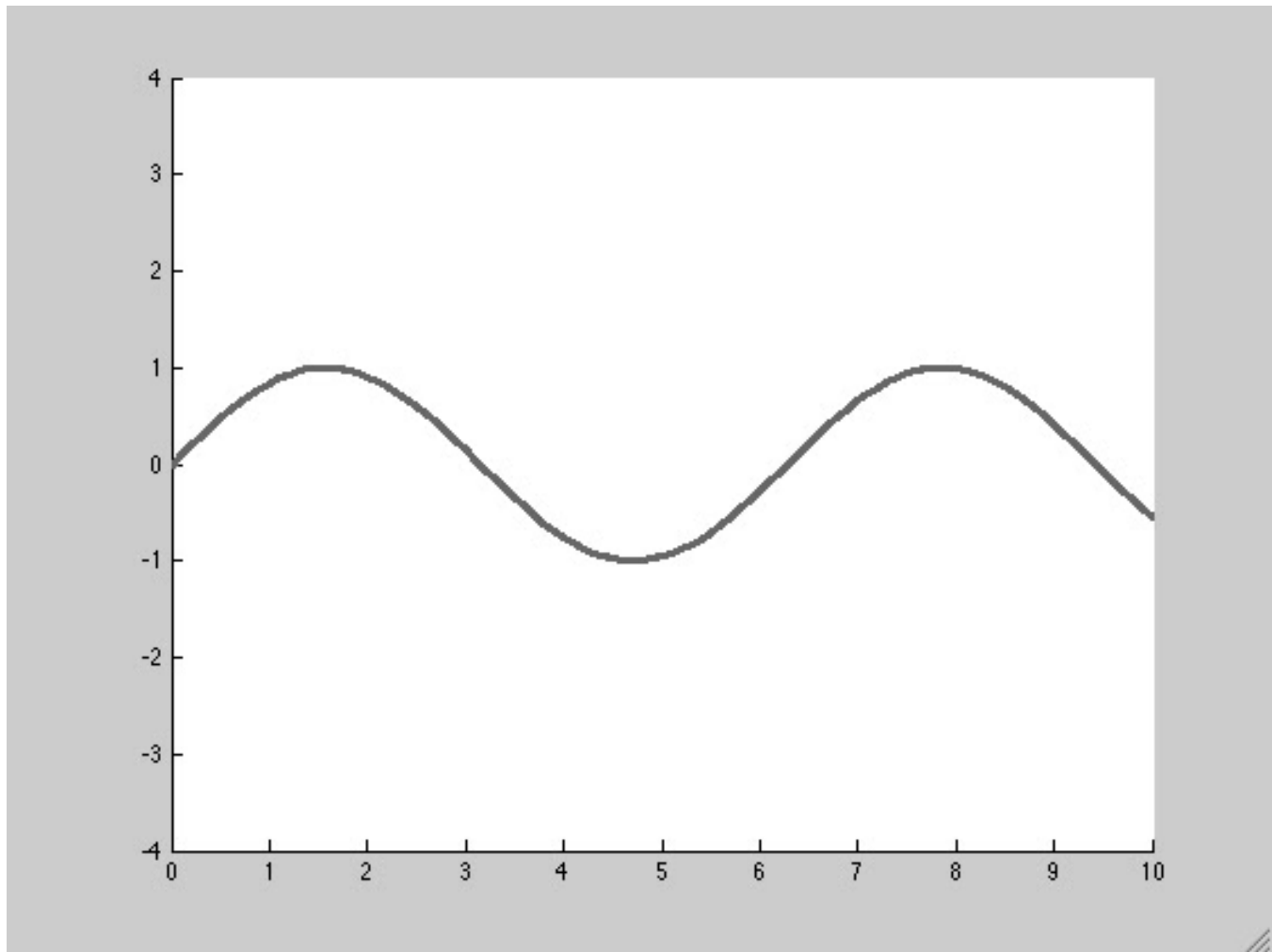
$$p(x_a | x_b) = N(\mu_{a|b}, \Sigma_{a|b})$$

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

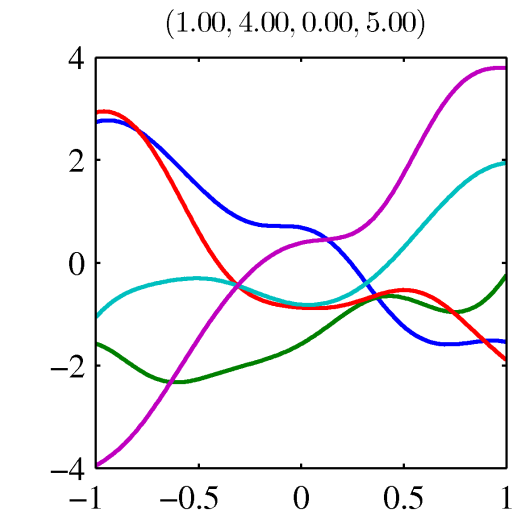
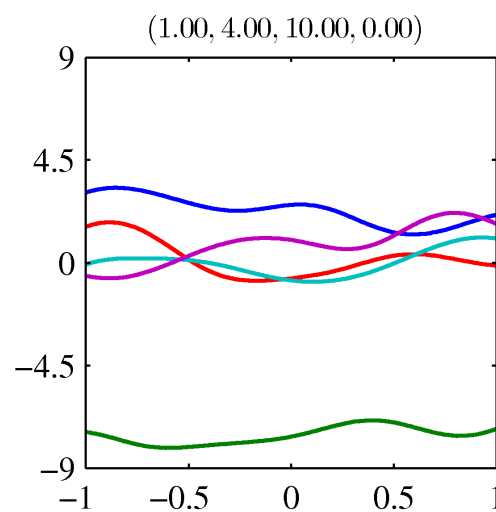
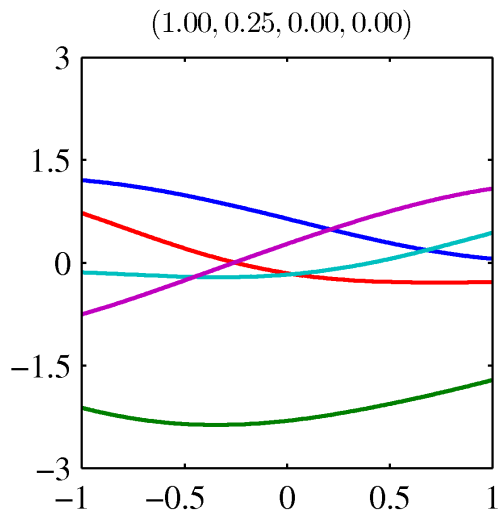
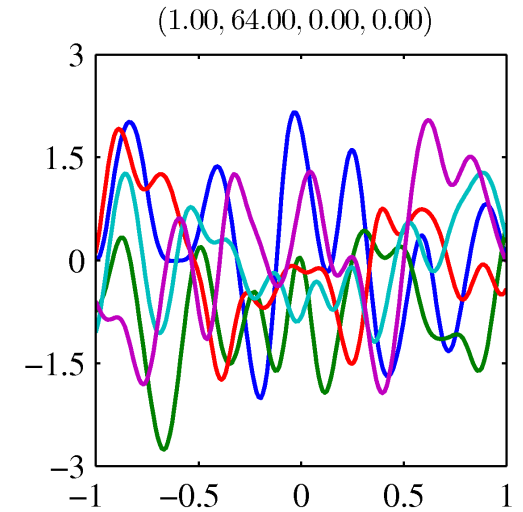
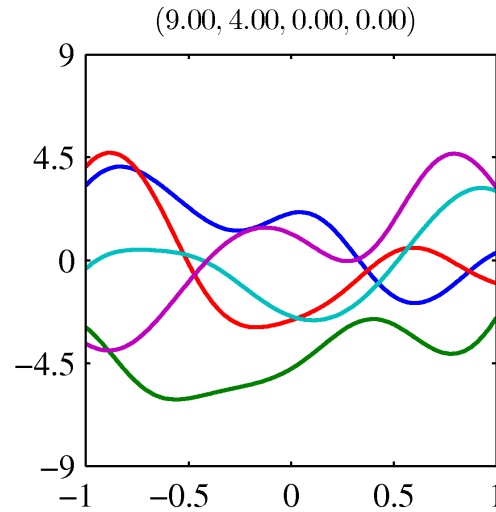
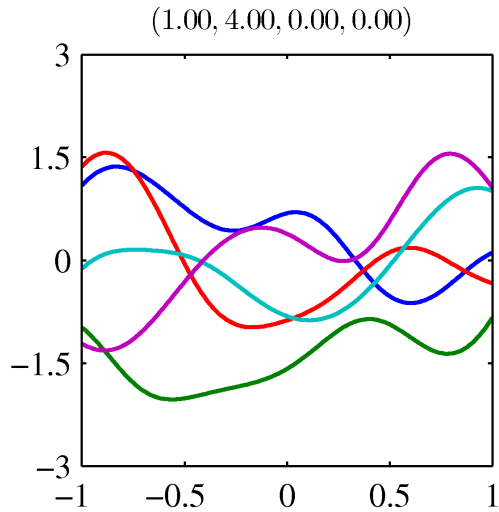
Recall conditional

GP Prediction



Hyperparameters

$$k(x, x') = \sigma_f^2 e^{-\frac{1}{2}(x-x')^T W (x-x')}$$



Hyperparameter Estimation

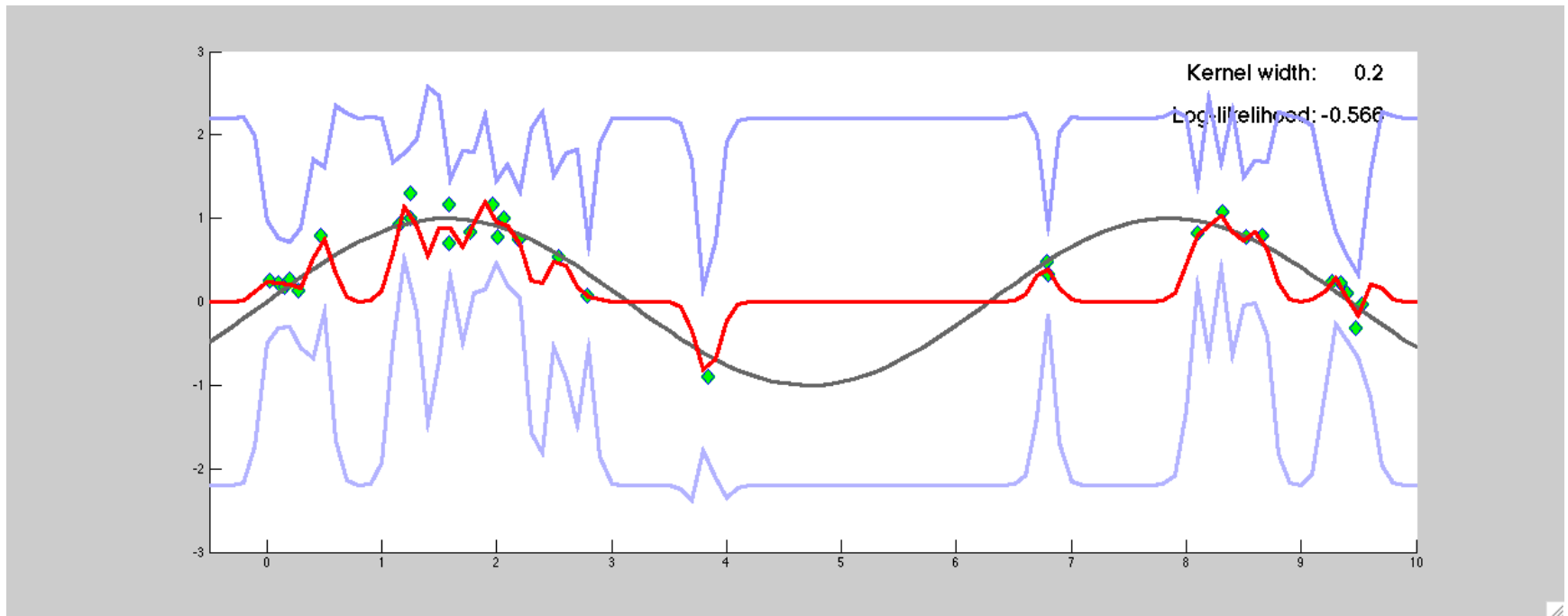
- Maximize data log likelihood:

$$\theta_* = \arg \max_{\theta} p(\mathbf{y} | \mathbf{X}, \theta)$$

$$\log p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log(\mathbf{K} + \sigma_n^2 \mathbf{I}) - \frac{n}{2} \log 2\pi$$

- Compute derivatives wrt. params $\theta = \langle \sigma_n^2, l, \sigma_f^2 \rangle$
- Optimize using conjugate gradient descent

Kernel Width



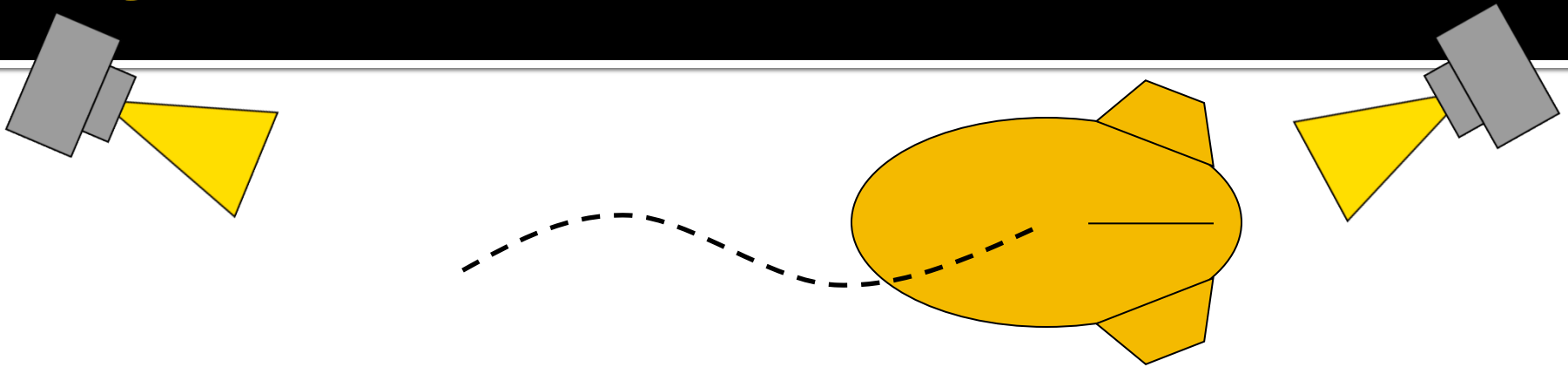
GP Optimization



- Learn hyperparameters via numerical methods
- Learn noise model at the same time

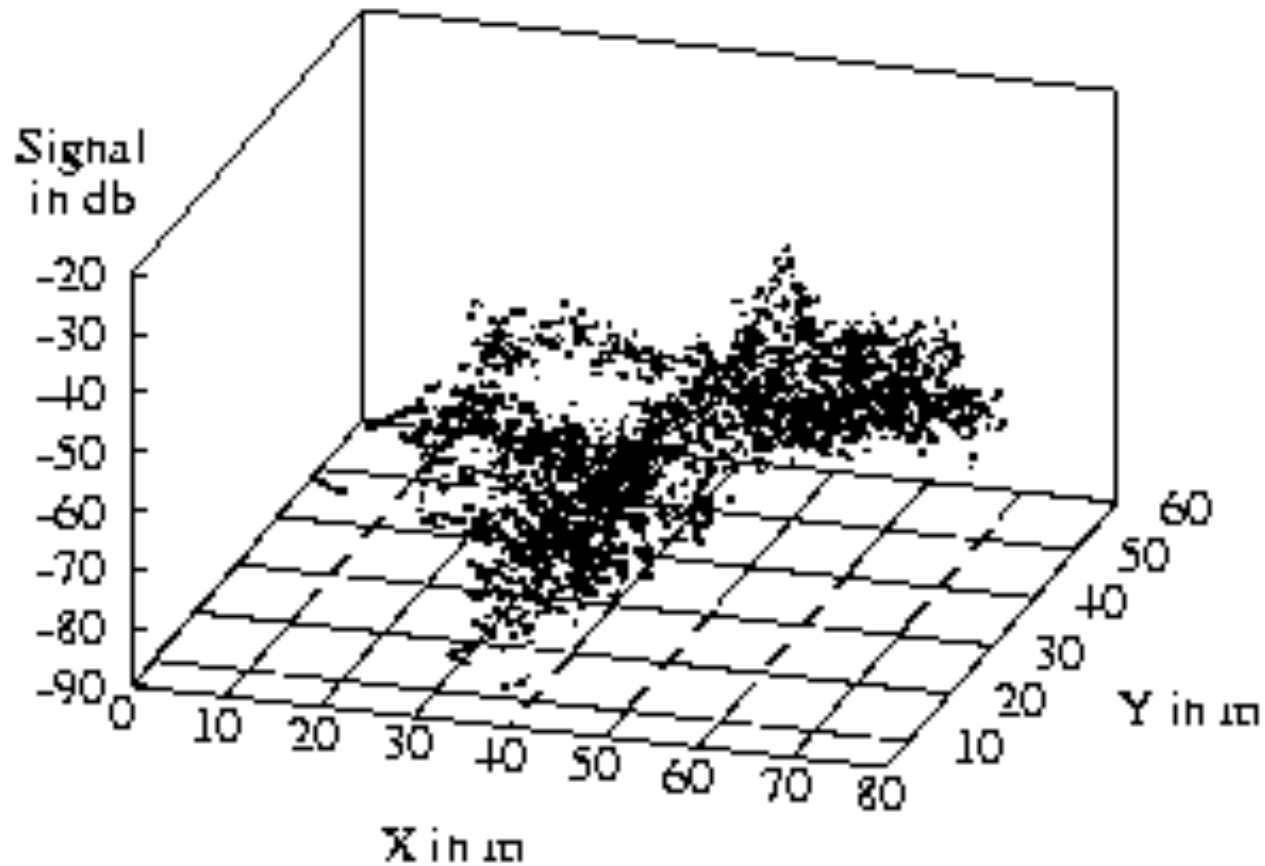
Bayesian Filtering with GPs

Dynamics and Observation Models

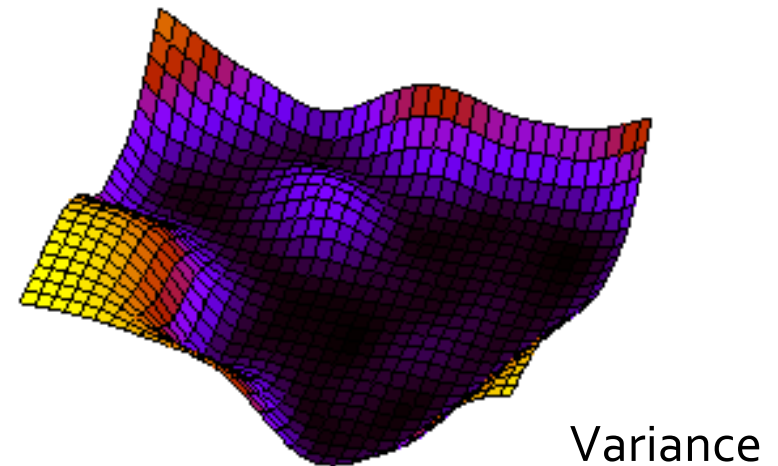
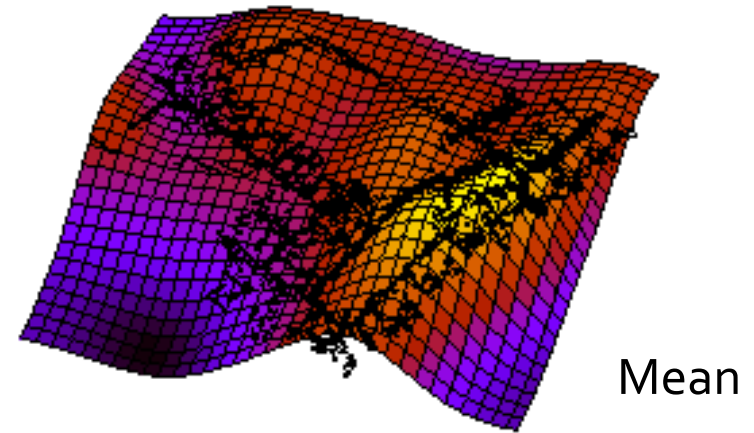
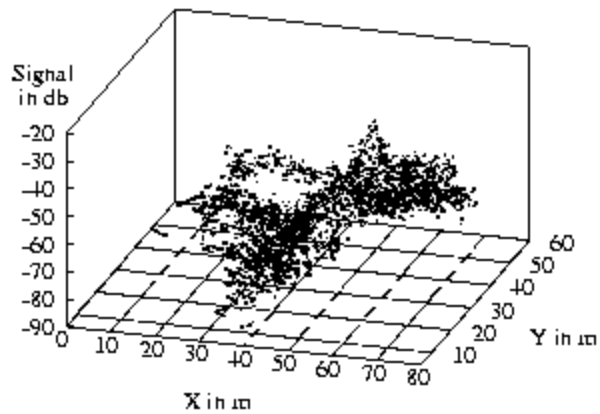


- System modeling
 - Learn behavior of the system given ground truth states
- Dynamics model
 - Discrete-time model of change over time
- Observation model
 - Mapping between states and observations
- Traditionally use parametric models
 - Derive system equations then learn parameters

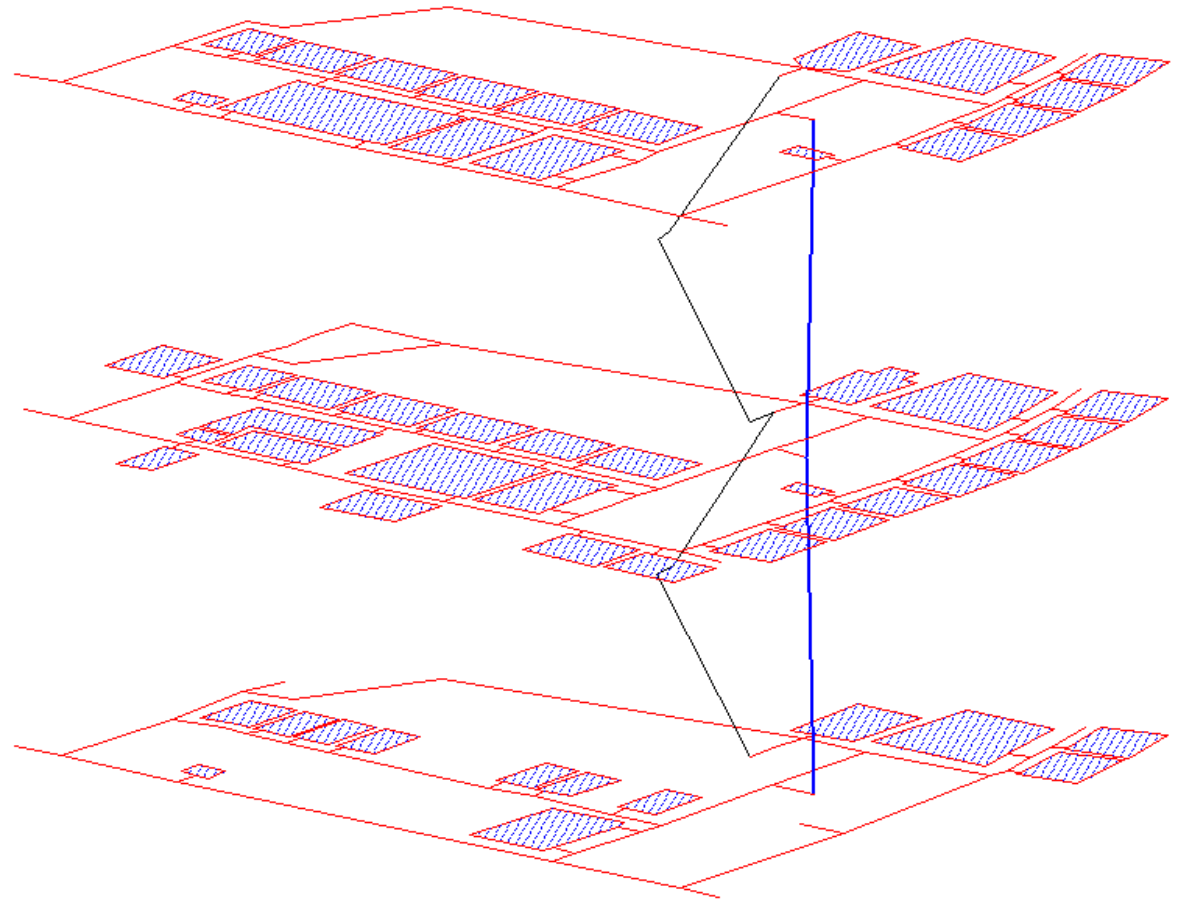
WiFi Sensor Model



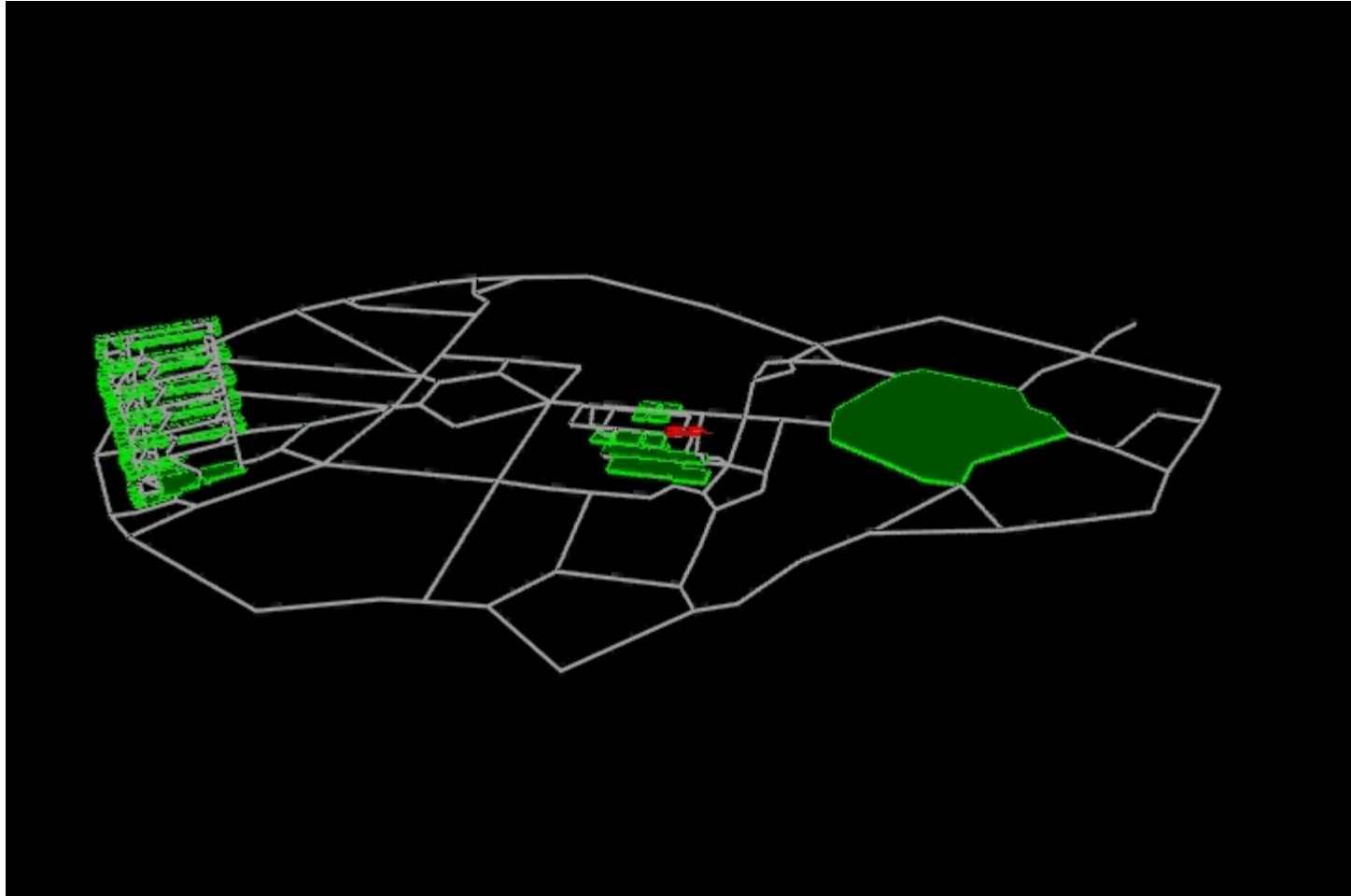
WiFi Sensor Model



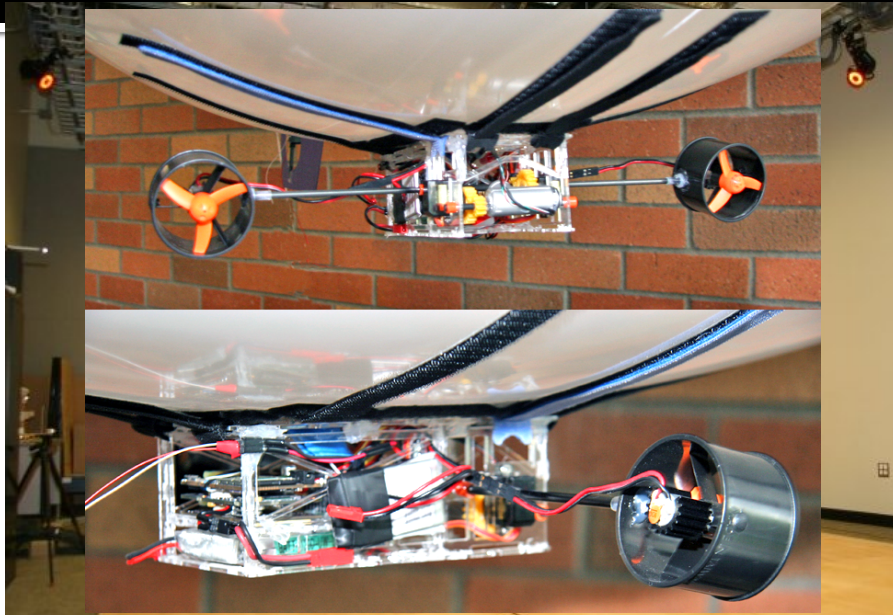
Mixed Representation



Tracking Example



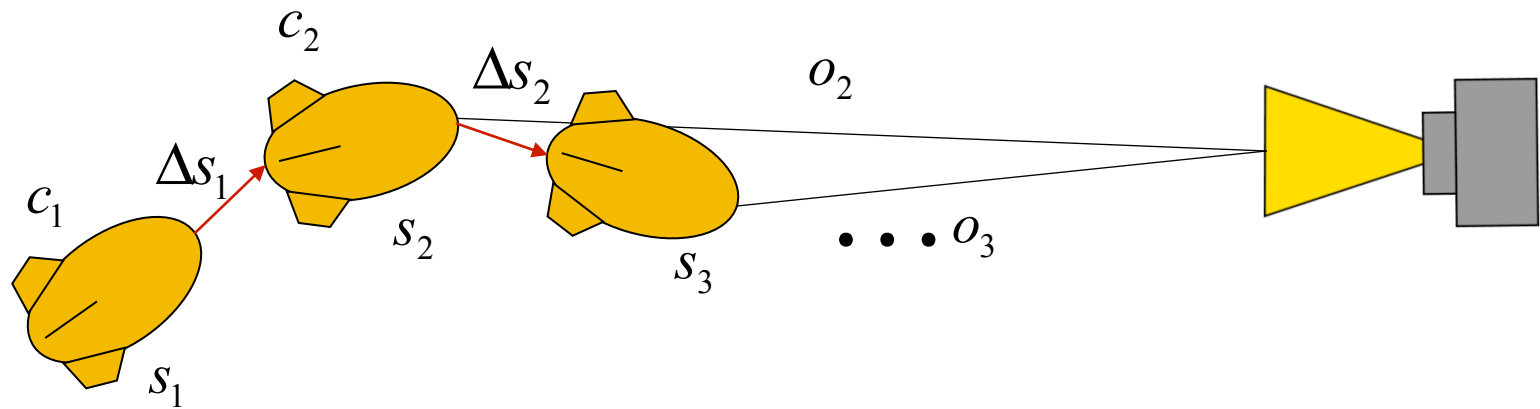
Blimp Test Platform



observation

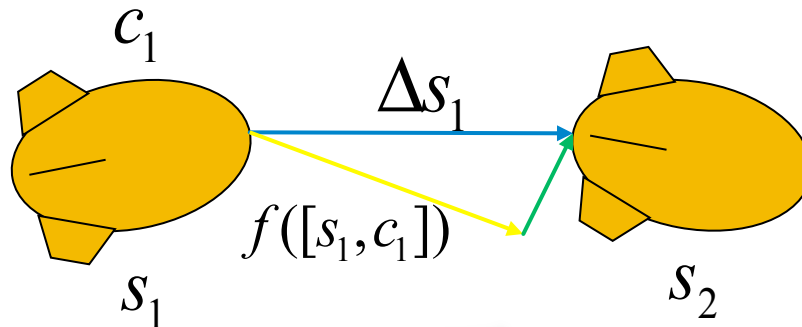
- System:
 - Commercial blimp envelope with custom gondola
 - XScale based computer with Bluetooth connectivity
 - Two main motors with tail motor (3D control)
- Observations:
 - Two cameras each operating at 1Hz
 - Extract ellipse using computer vision (5D observations)
- Ground truth obtained via VICON motion capture system

Learning GP Dynamics and Observation Models



- Use ground truth state to extract:
 - Dynamics data
$$D_s = \langle [s_1, c_1], \Delta s_1 \rangle, \langle [s_2, c_2], \Delta s_2 \rangle \dots$$
 - Observation data
$$D_o = \langle s_2, o_2 \rangle, \langle s_3, o_3 \rangle$$
- Learn models using Gaussian process regression
 - Learn process noise inherent in system

Learning Enhanced-GP Models



- Combine GP model with parametric model

$$D_x = \langle [s_1, c_1], \Delta s_1 - f([s_1, c_1]) \rangle$$

- Advantages

- Captures aspects of system not considered by parametric model
- Learns noise model in same way as GP-only models
- Higher accuracy for same amount of training data

GP Modeling Accuracy

Dynamic model error

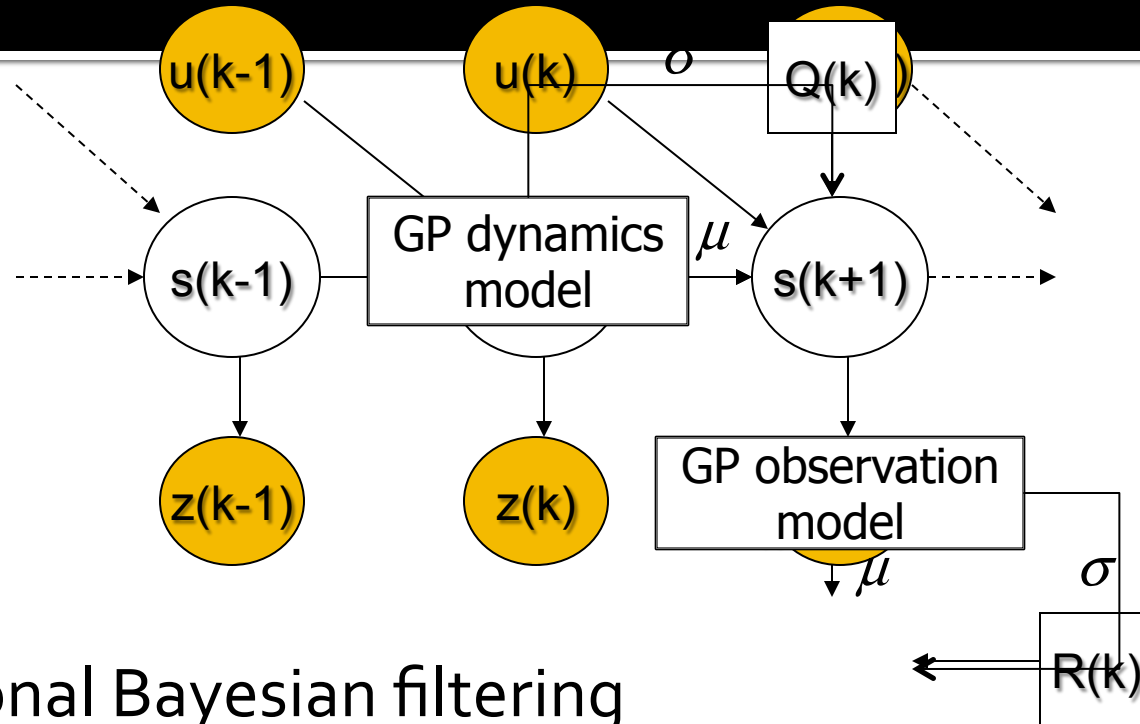
Propagation method	pos(mm)	rot(deg)	vel(mm/s)	rotvel(deg/s)
Param	3.3	0.5	14.6	1.5
GPonly	1.8	0.2	9.8	1.1
EGP	1.6	0.2	9.6	1.3

Observation model error

Modeling method	pos(pix)	Major axis(pix)	Minor axis(pix)	Theta(deg)
Param	7.1	2.9	5.7	9.2
GPonly	4.7	3.2	1.9	9.1
EGP	3.9	2.4	1.9	9.4

- 1800 training points, mean error over 900 test points
- For dynamic model, 0.25 sec predictions

GP-BayesFilters

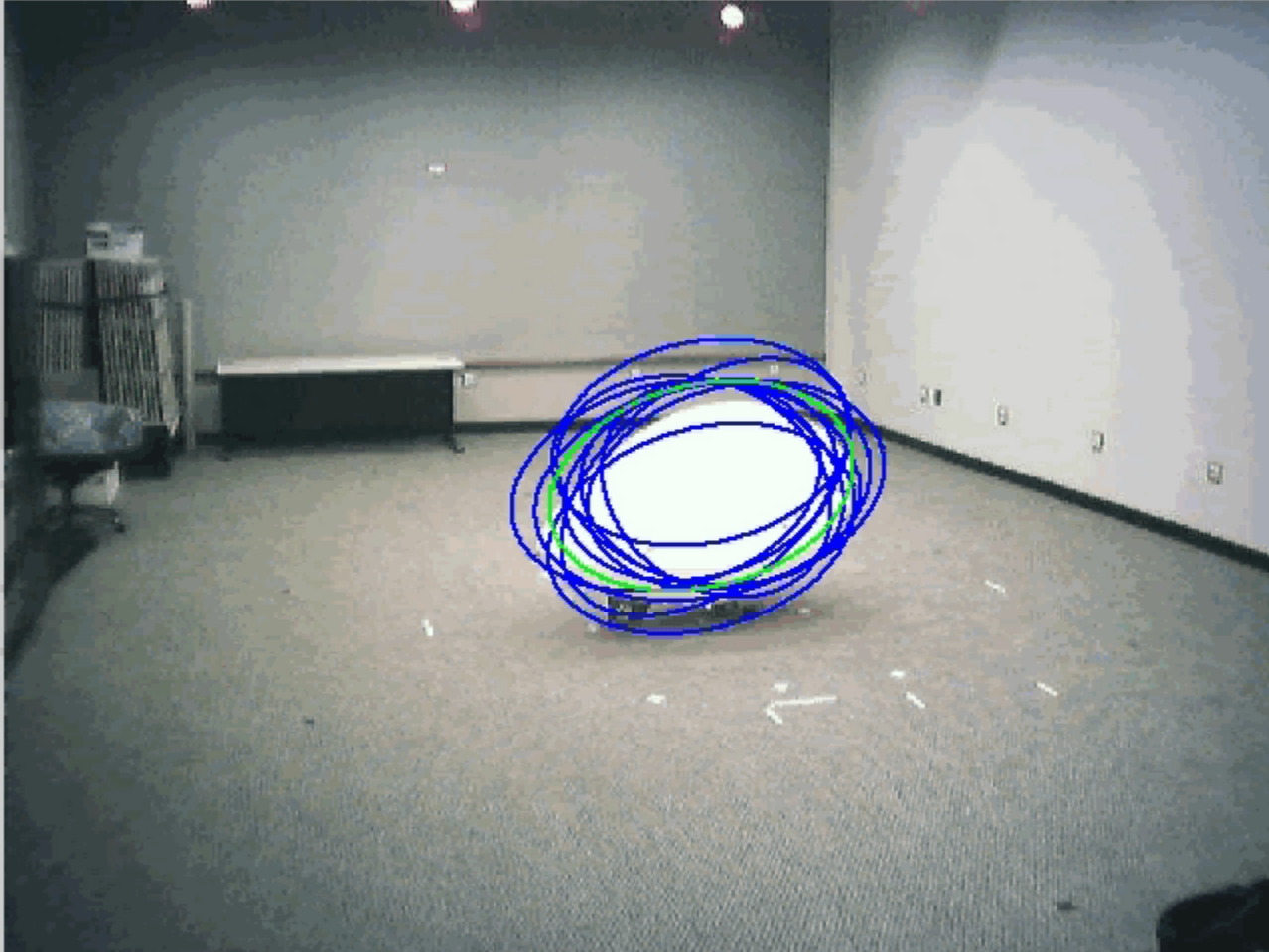


- Traditional Bayesian filtering
 - Parametric dynamics and observation models
- GP-BayesFilters
 - GP dynamics and observation models
 - Noise derived from GP prediction uncertainty
 - Can be integrated into to Bayes filters: EKF, UKF, PF, ADF

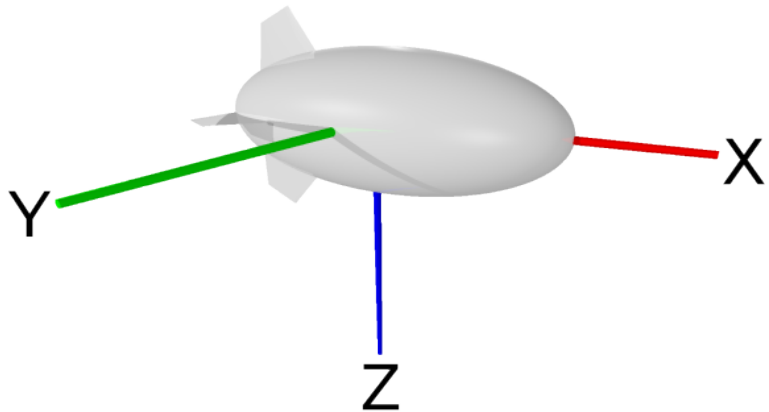
GP-BayesFilters

- **Learn GP:**
 - Input: Sequence of **ground truth states** along with controls and observations: $\langle s, u, z \rangle$
 - Learn GPs for dynamics and observation models
- **Filters**
 - **Particle filter:** sample from dynamics GP, weigh by Gaussian GP observation function
 - **EKF:** GP for mean state, GP derivative for linearization
 - **UKF:** GP for sigma points

GP-UKF Tracking Example



Non-linear Parametric Model

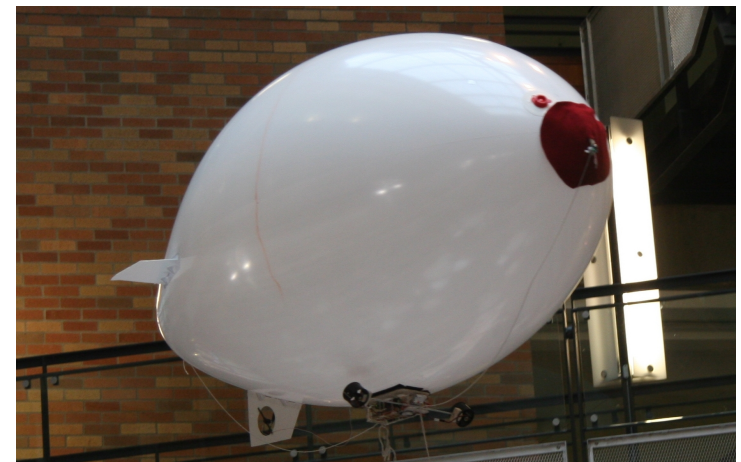

$$\dot{s} = \frac{d}{dt} \begin{bmatrix} p \\ \xi \\ v \\ \omega \end{bmatrix} = \begin{bmatrix} R_b^e v \\ H(\xi) \\ M^{-1}(\sum Forces - \omega^* M v) \\ J^{-1}(\sum Torques - \omega^* J \omega) \end{bmatrix}$$

- 12-D state=[pos,rot,transvel,rotvel]
- Describes evolution of state as ODE
- Forces / torques considered: buoyancy, gravity, drag, thrust
- 16 parameters are learned by optimization on ground truth motion capture data

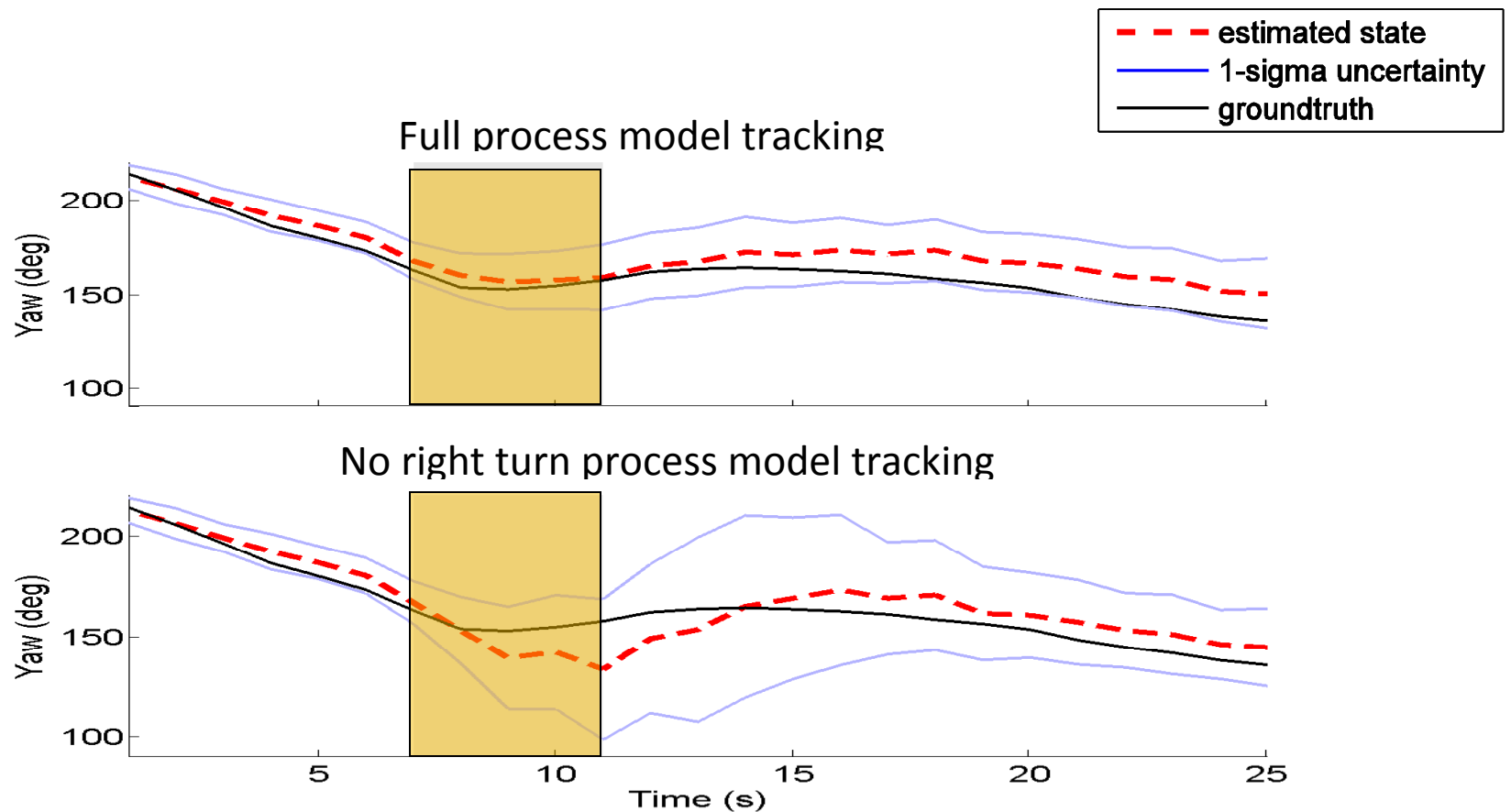
Blimp Results

Tracking algorithm	pos(mm)	rot(deg)	vel(mm/s)	rotvel(deg/s)	MLL	time(sec)
GP-PF	91 \pm 7	6.4 \pm 1.6	52 \pm 3.7	5.0 \pm .2	9.4 \pm 1.9	449.4 \pm 21
GP-EKF	93 \pm 1	5.2 \pm .1	52 \pm .5	4.6 \pm .1	13.0 \pm .2	.29 \pm .1
GP-UKF	89 \pm 1	4.7 \pm .2	50 \pm .4	4.5 \pm .1	14.9 \pm .5	1.28 \pm .3
ParaPF	115 \pm 5	7.9 \pm .1	64 \pm 1.2	7.6 \pm .1	-4.5 \pm 4.2	30.7 \pm 5.8
ParaEKF	112 \pm 4	8.0 \pm .2	65 \pm 2	7.5 \pm .2	8.4 \pm 1	.21 \pm .1
ParaUKF	111 \pm 4	7.9 \pm .1	64 \pm 1	7.6 \pm .1	10.1 \pm 1	.33 \pm .1

- Blimp tracking using multiple cameras
- Ground truth obtained via Vicon motion tracking system
- Average tracking error
- Trajectory ~12 min long
- 0.5 sec timesteps



Dealing With Training Data Sparsity



- Training data for right turns removed

Going Latent

- Sometimes ground truth states are not or only partially available.
- Instead of optimizing over hyperparameters only

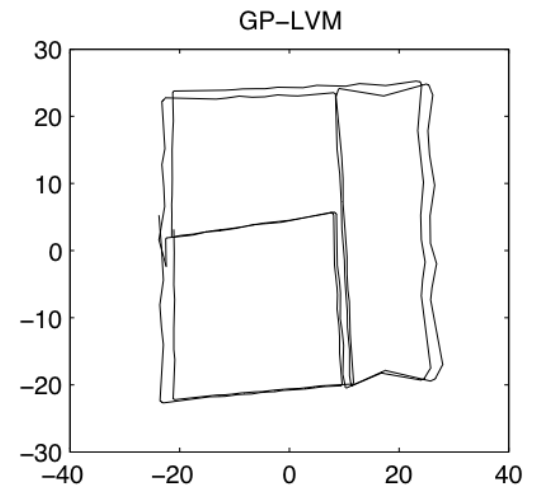
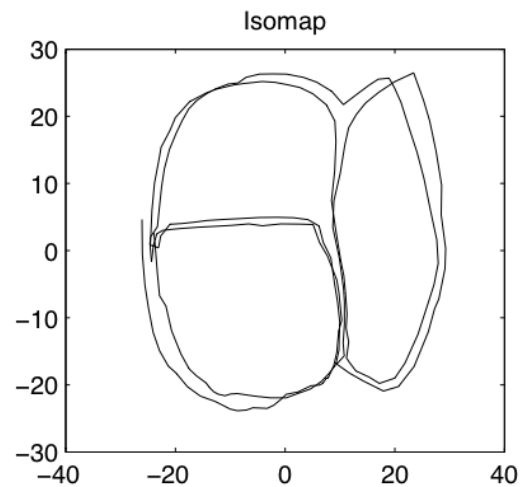
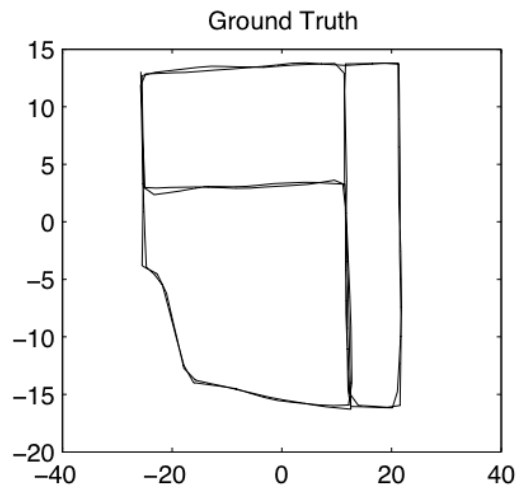
$$\theta_* = \arg \max_{\theta} p(\mathbf{y} | \mathbf{X}, \theta)$$
 optimize over latent states X as well

$$\langle \mathbf{X}_*, \theta_* \rangle = \arg \max_{\mathbf{X}, \theta} p(\mathbf{y} | \mathbf{X}, \theta)$$
- **GPLVM**: non-linear probabilistic dimensionality reduction

WiFi-SLAM:

[Ferris-Fox-Lawrence: IJCAI-07]

Mapping without Ground Truth



GPDM: Latent Variable Models for Dynamical Systems

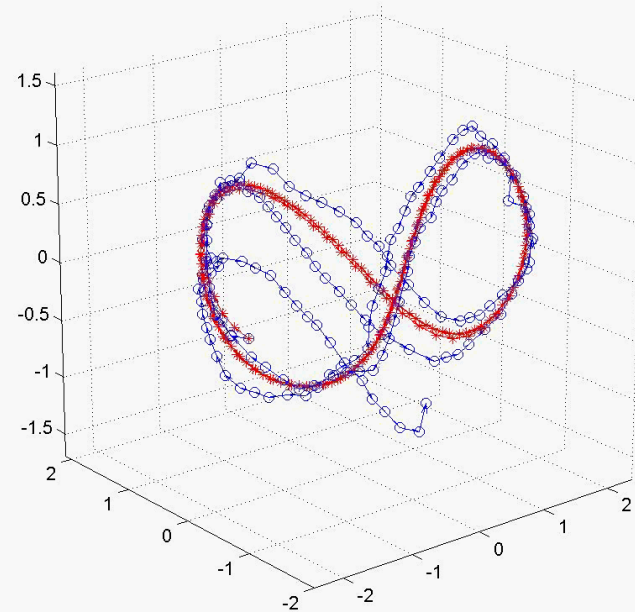
- GPLVM with additional constraints on latent variables to model dynamical system

$$p(\mathbf{y}, \mathbf{X}, \theta) = p(\mathbf{y} | \mathbf{X}, \theta) p(\mathbf{X} | \theta) p(\theta)$$

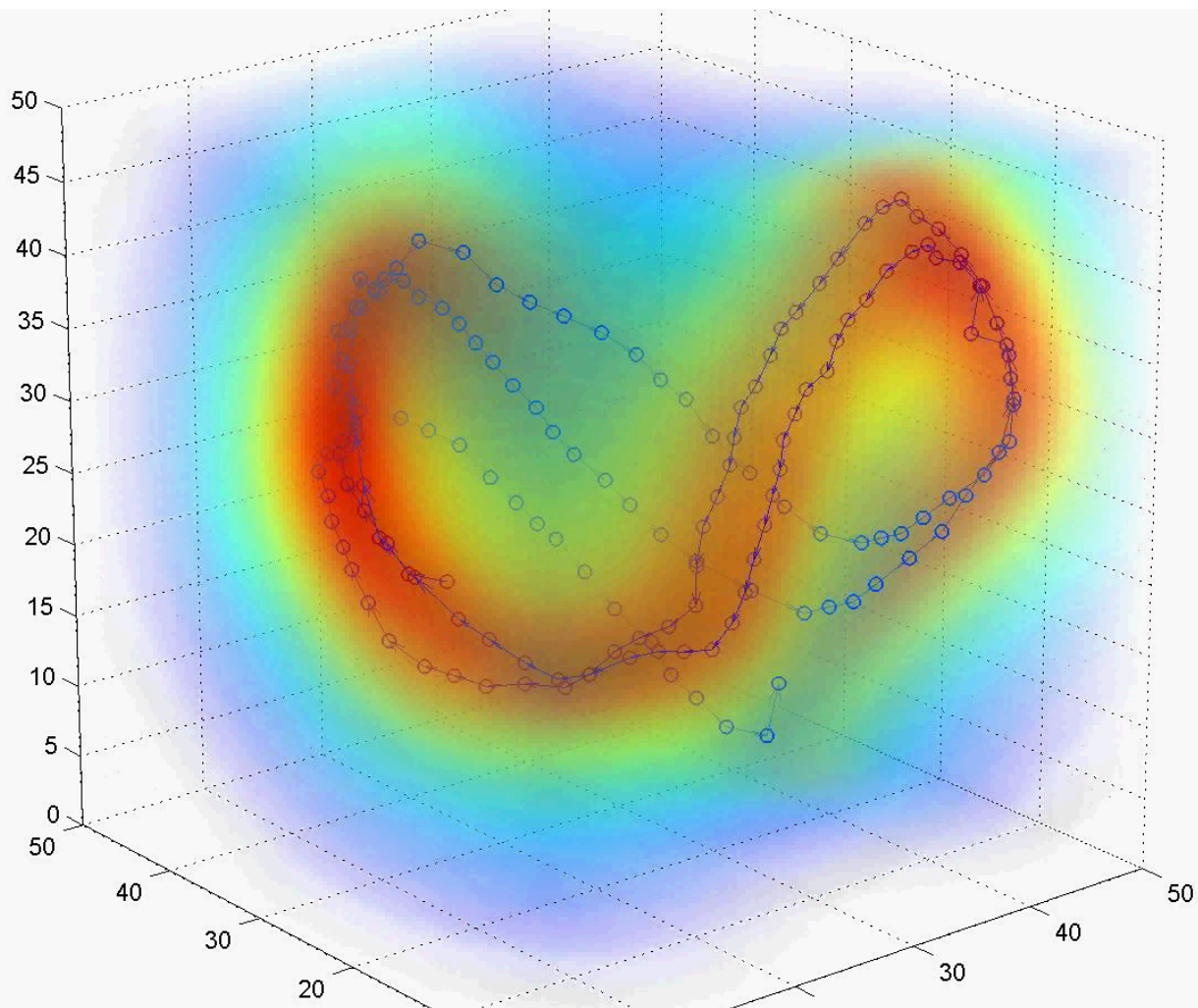
$$\langle \mathbf{X}_*, \theta_* \rangle = \arg \max_{\mathbf{X}, \theta} p(\mathbf{y}, \mathbf{X}, \theta)$$

- Dynamics are modeled via another GP

GPDM Application



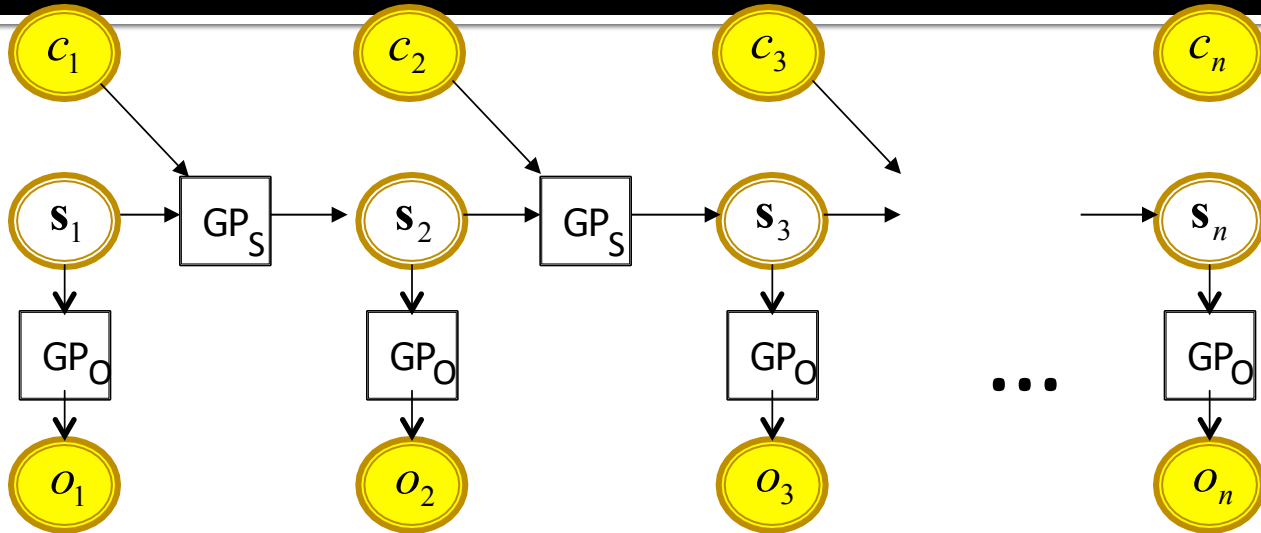
GPDM Volume



GPBF-Learning

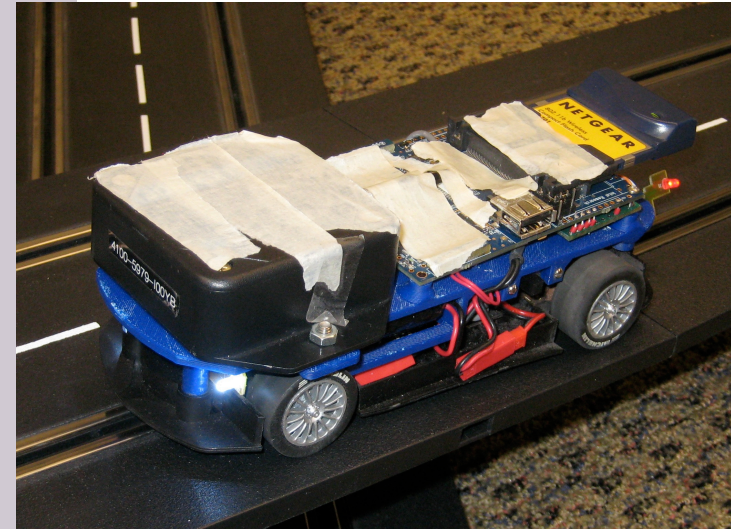
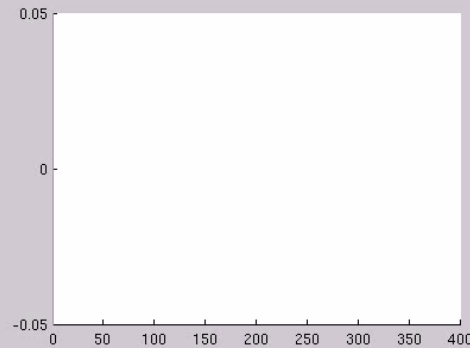
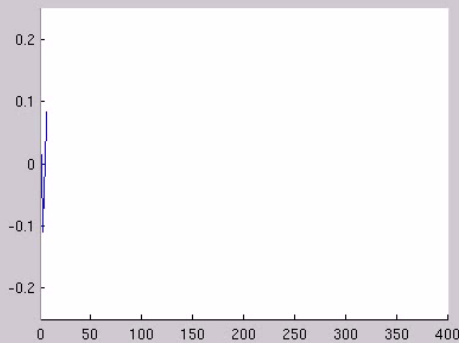
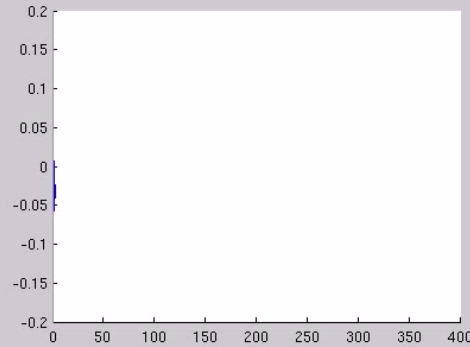
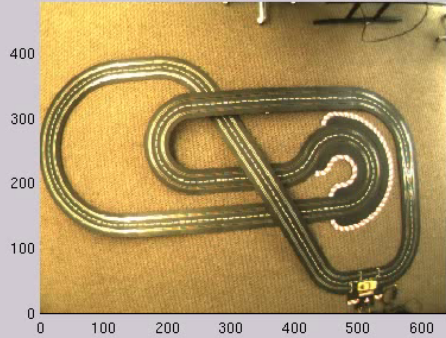
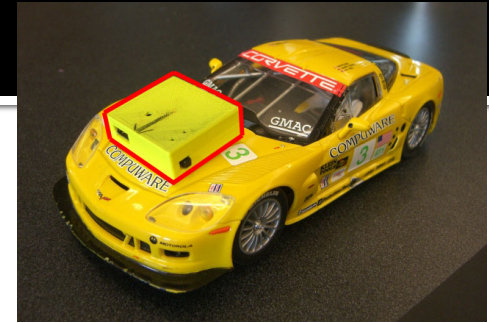
- Extend GPDMs to
 - incorporate control
 - incorporate sparse labels on latent states
- Steps:
 - Learn GPLVM
 - Extract GP-BayesFilter for tracking

GPBF-Learn



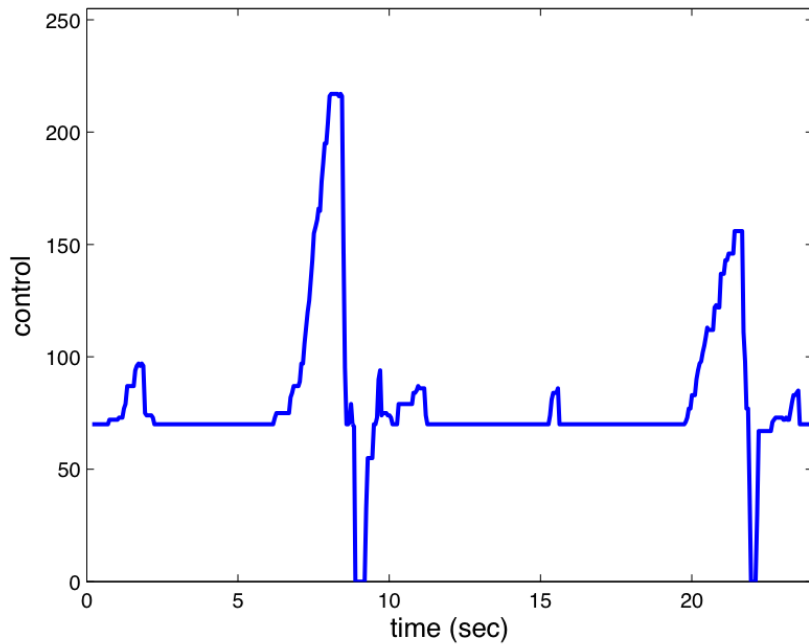
- GPBF-Learn follows same structure as Bayes filter
- Learns latent states and GPs in one optimization
 - GP_S for system dynamics
 - GP_O for observation model
- Can perform system identification

Evaluation System: Slotcar

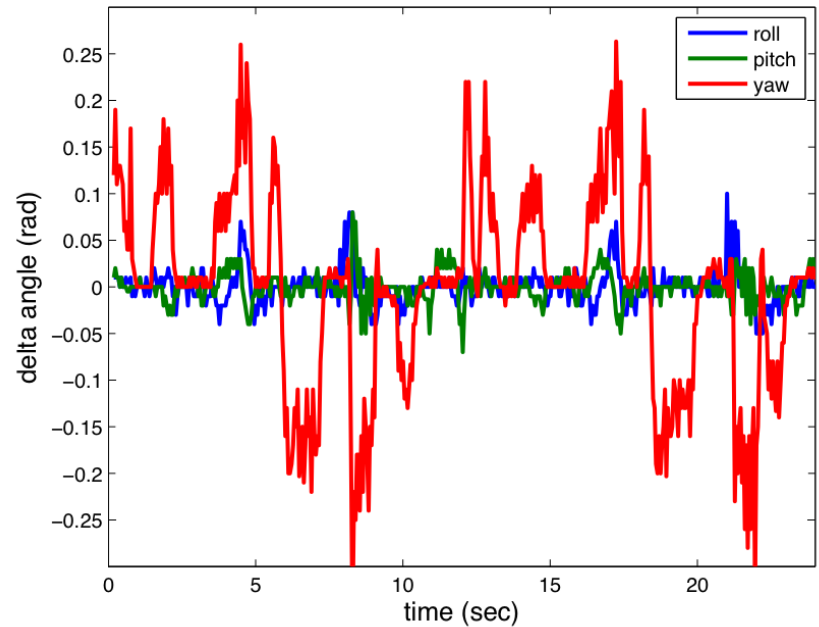


- Track contains banked curves, elevation changes
- Custom IMU with gyros and accelerometers
- Observations very noisy, perceptual aliasing

Data

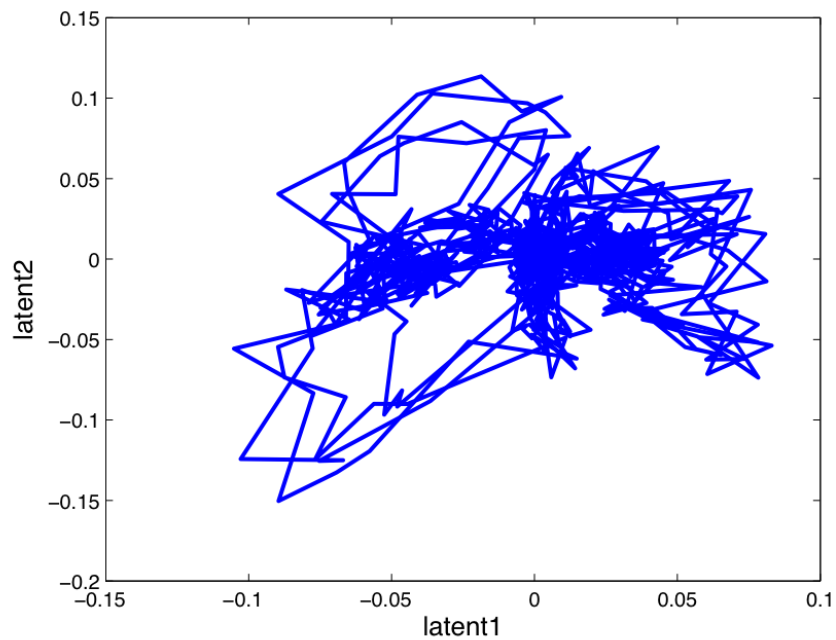


Controls

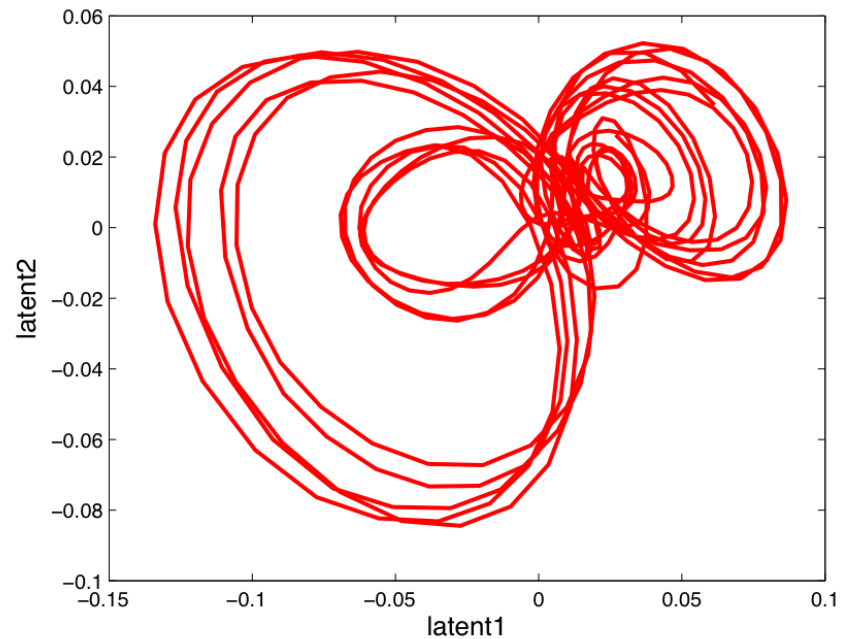


IMU: angle change rate

Learned 2D Embedding

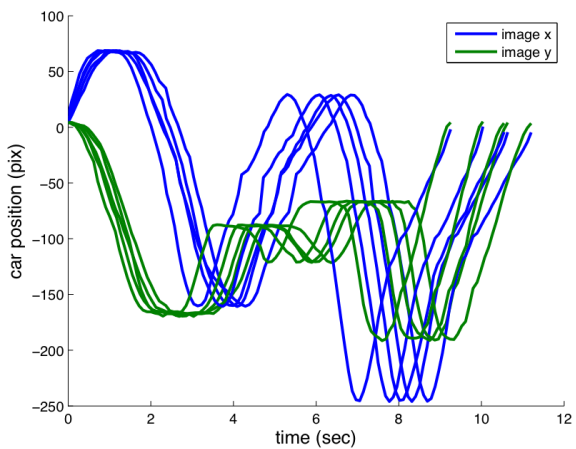


N4SID

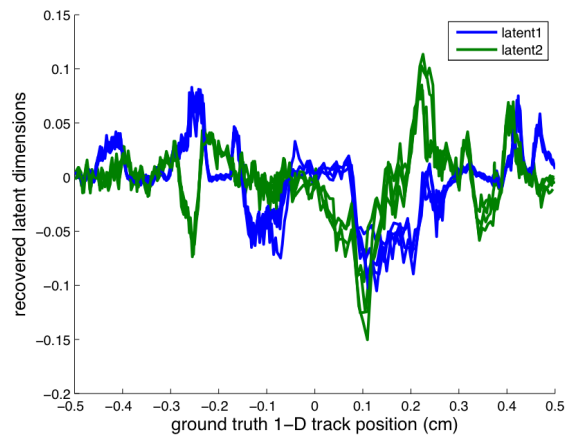


GPBF-Learn

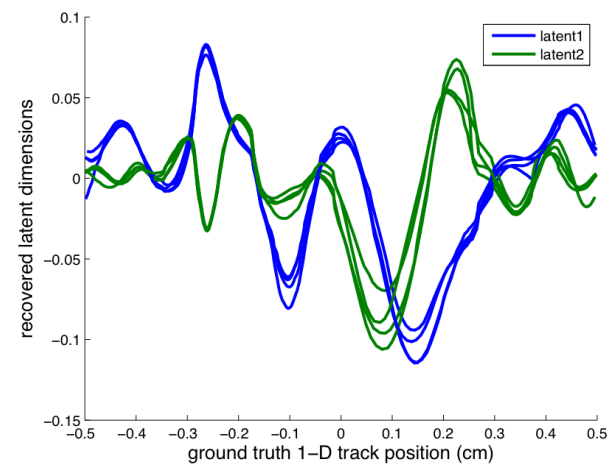
Time Alignment



Raw

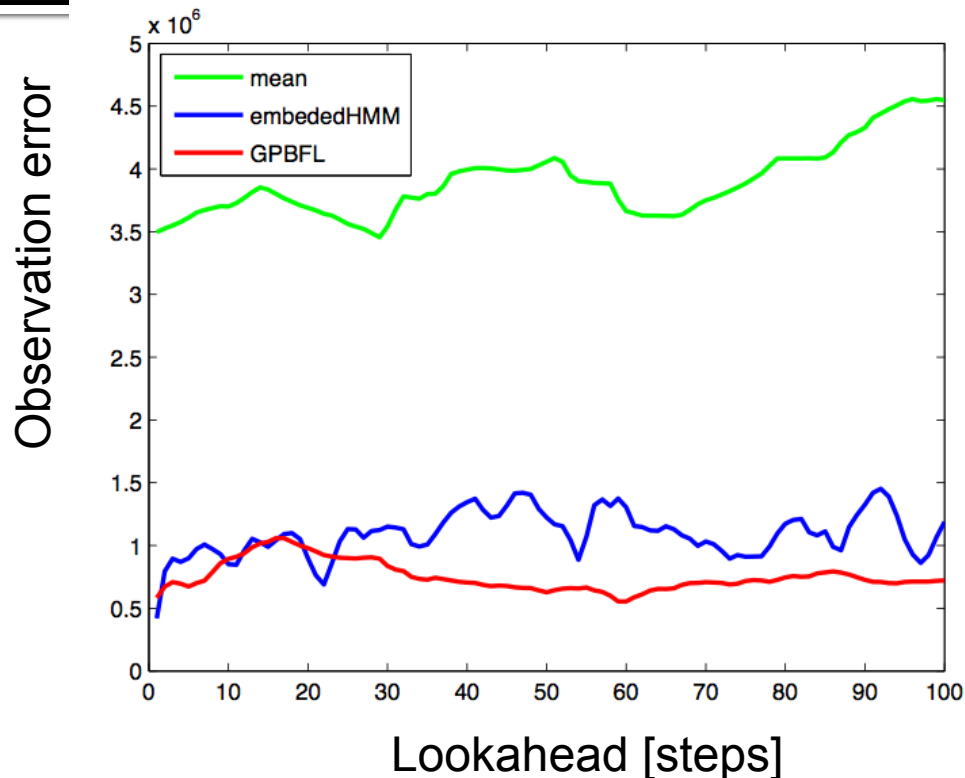


N4SID



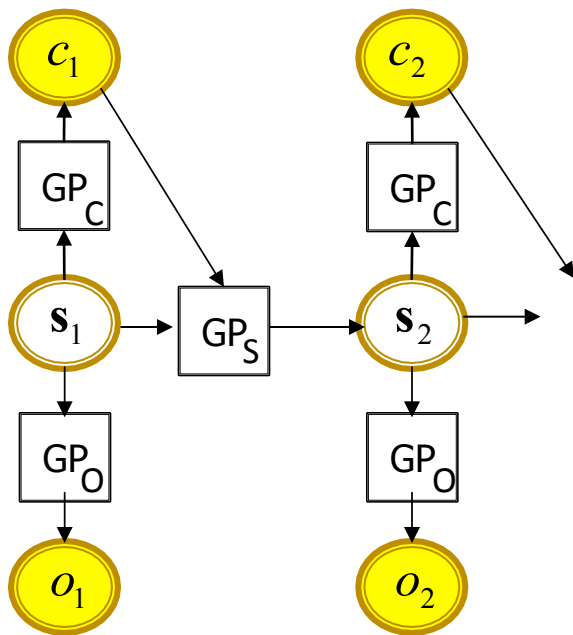
GPBF-Learn

GPBFL and HSE-HMMs



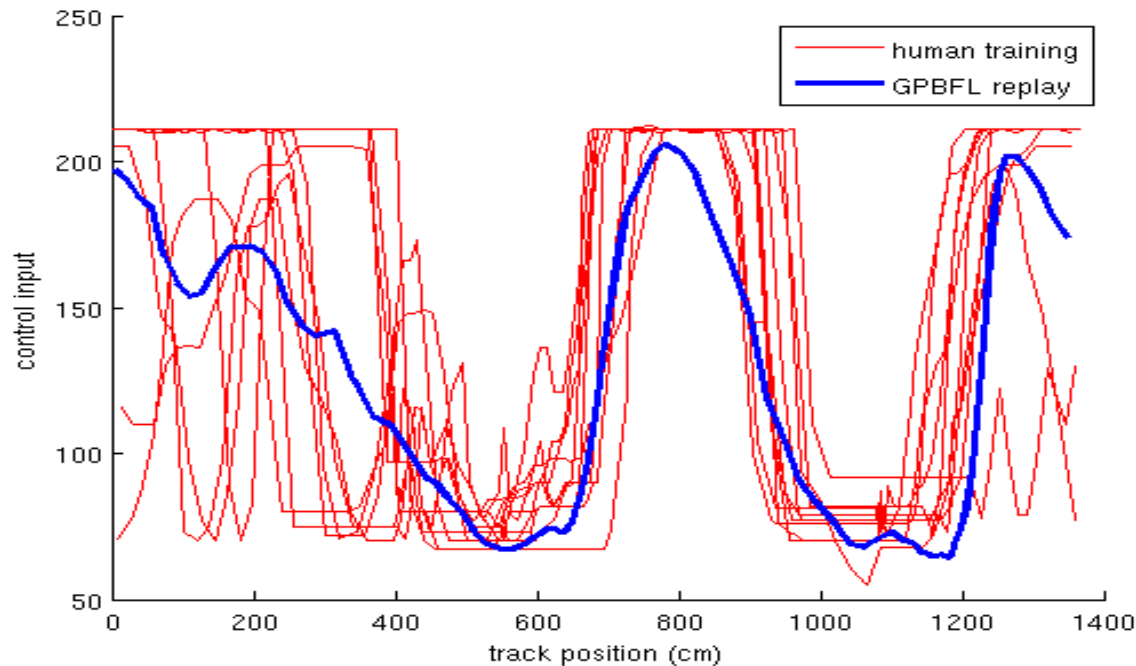
- Superior to Kernel CCA [Kawahara et al., NIPS-07]
- Latent space dimensionality: GPBFL 3D, HSE-HMM 20D
- HSE-HMM much more efficient!

Learning by Demonstration with GPBF-Learn



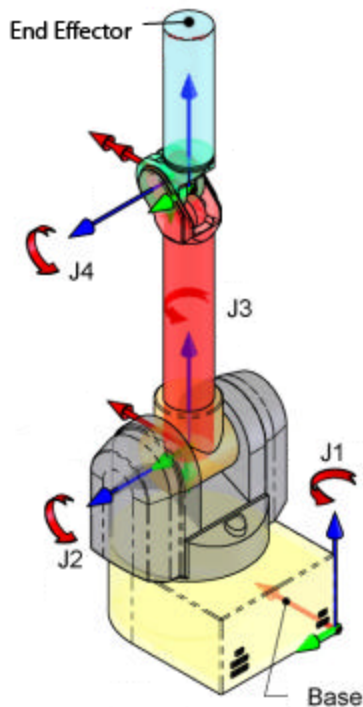
- Learning strategy
 - Learn latent states using GPBF-Learn
 - Learn mapping from state to control
- Replay strategy
 - track state using GP-BayesFilter
 - use control given by control model
- Requires only observations and user demonstrated control

Simple Control Experiment



- Learn control for slot car given human demonstrations
- Replay interpolates human examples
- Fails for more complicated systems

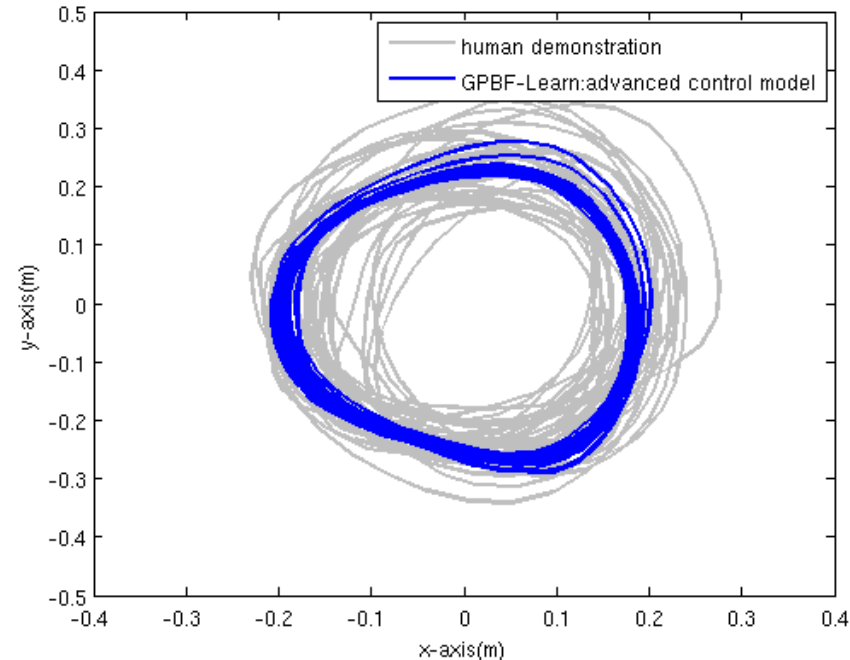
WAM Test Platform



- System:
- Barrett Whole Arm Manipulator
 - Four joints/degrees of freedom
 - 4D control (change in joint angles)
 - Significant control noise
- Observations:
 - 3D position of end effector
- User demonstration:
 - Manipulate to trace out circular trajectory for end effector

Advanced Control Experiment

Top down view of end effector position



- Learn 3D latent states for system
- Time-based replay fails
- Advanced control model achieves proper replay

Beyond Plain GP Regression

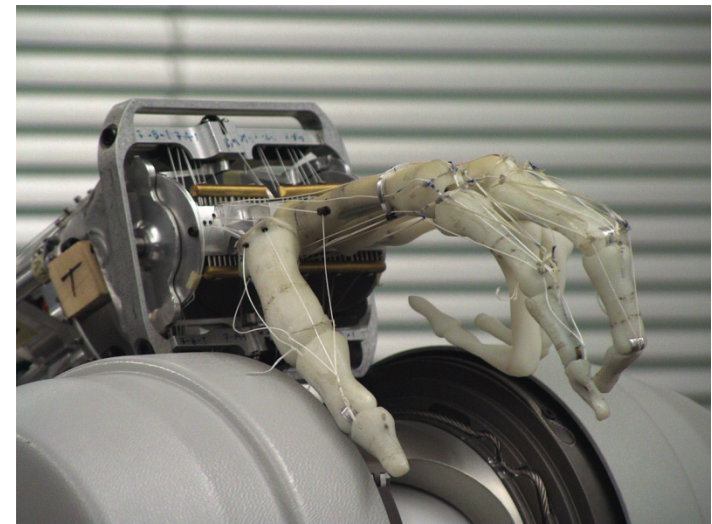
- **Heteroscedastic** (state dependent) noise
[Quoc-etal: ICML-05; Kersting-etal: IMCL-07]
- Efficiency: **Sparse** GPs
[Snelson and Ghahramani: NIPS-06; Smola and Bartlett: NIPS-01]
- **Discrete** state components: GP classification
[Plagemann etal: IJCAI-07]

Other Issues

- Heteroscedastic (state dependent) noise
- Non-stationary GPs
- Coupled outputs
- Sparse GPs
 - Online: Decide whether or not to accept new point
 - Remove points
 - Optimize small set of points
- Classification
 - Laplace approximation
 - No closed-form solution, sampling

Summary

- GPs provide **flexible modeling framework**
- Take **data noise and uncertainty due to data sparsity** into account
- Combination with parametric models increases accuracy and reduces need for training data
- Computational complexity is a key problem



Some References

- Website: <http://www.gaussianprocess.org/>
- GP book: <http://www.gaussianprocess.org/gpml/>
- GPLVM: <http://www.cs.man.ac.uk/~neill/gplvm/>
- GPDM: <http://www.dgp.toronto.edu/~jmwang/gpdm/>
- Bishop book:
<http://research.microsoft.com/en-us/um/people/cmbishop/prml/>