



UNIVERSITÄT DARMSTADT

Approximate Input Inference as Stochastic Optimal Control

Joe Watson

Intelligent Autonomous Systems TU Darmstadt watson@ias.tu-darmstadt.de ♥ @JoeMWatson, @ias_tudarmstadt



Figure 1: Object Manipulation



Figure 2: Full-body Locomotion





Figure 1: Object Manipulation

Figure 2: Full-body Locomotion

$$\min_{\mathbf{u}_{0:T-1}} C_T(\boldsymbol{x}_T) + \sum_{t=0}^{T-1} C(\boldsymbol{x}_t, \mathbf{u}_t)$$

s.t. $\boldsymbol{x}_{t+1} = \mathbf{f}(\boldsymbol{x}_t, \mathbf{u}_t)$





Figure 1: Object Manipulation

Figure 2: Full-body Locomotion

$$\min_{\mathbf{u}_{0:T-1}} C_T(\boldsymbol{x}_T) + \sum_{t=0}^{T-1} C(\boldsymbol{x}_t, \mathbf{u}_t)$$

s.t. $\boldsymbol{x}_{t+1} = \mathbf{f}(\boldsymbol{x}_t, \mathbf{u}_t)$





Figure 1: Object Manipulation

Figure 2: Full-body Locomotion

$$\min_{\mathbf{u}_{0:T\cdot 1}} C_T(oldsymbol{x}_T) + \sum_{t=0}^{T-1} C(oldsymbol{x}_t, \mathbf{u}_t)$$

s.t. $oldsymbol{x}_{t+1} \sim \mathbf{f}(oldsymbol{x}_t, \mathbf{u}_t)$

Many successful Model-based Reinforcement Learning (MBRL) formulations combine probabilistic model learning with trajectory optimization:

Algorithm	Probabilistic Model	Optimizer	Weakness
PILCO ¹	GP	Backprop-through-time	Unstable gradients
AGP-iLQR ²	GP	iLQR	Doesn't use uncertainty
Guided Policy Search ³	GMM	Maximum Entropy iLQG	Heavy regularization
$PETS^4$	NN Ensemble	CEM	Computational cost

⁴Marc Deisenroth and Carl E Rasmussen. "PILCO: A model-based and data-efficient approach to policy search". In: Proceedings of the 28th International Conference on machine learning (ICML-11). 2011.

⁴Joschka Boedecker et al. "Approximate real-time optimal control based on sparse gaussian process models". In: 2014 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL). IEEE. 2014.

⁴Sergey Levine. "Motor skill learning with local trajectory methods". PhD thesis. Stanford University, 2014.

⁴Kurtland Chua et al. "Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models". In: Advances in Neural Information Processing Systems 31. Probabilistic Numerics for Trajectory Optimization:

Probabilistic Numerics⁵

Numerical methods tackled as probabilistic inference. The benefits are

- Uncertainty quantification
- Principled regularization
- Inclusion of prior knowledge through Priors

⁵Philipp Hennig, Michael A Osborne, and Mark Girolami. "Probabilistic numerics and uncertainty in computations". In: Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences (2015).

Input Estimation

$$\max_{\mathbf{u}_{0:T-1},\boldsymbol{\theta}} P(\mathbf{u}_{0:T-1}, \boldsymbol{x}_{0:T}, \mathbf{y}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1},\boldsymbol{\theta}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} \underbrace{P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})}_{\text{Dynamics}} \prod_{t=0}^{T-1} \underbrace{P(\mathbf{y}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})}_{\text{Measurement}}$$

Input Estimation

$$\max_{\mathbf{u}_{0:T-1},\boldsymbol{\theta}} P(\mathbf{u}_{0:T-1}, \boldsymbol{x}_{0:T}, \mathbf{z}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1},\boldsymbol{\theta}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} \underbrace{P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})}_{\text{Dynamics}} \prod_{t=0}^{T-1} \underbrace{P(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})}_{\text{Control Cost}}$$

$$\max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0:T}, \mathbf{u}_{0:T-1}, \boldsymbol{z}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta}) \prod_{t=0}^{T-1} P(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})$$

- Applies to finite horizon optimal control problems, moving from x_0 to x_g over horizon T
- **z** represents the 'optimization states' and their desired trajectory (i.e. $\mathbf{z}_t = [\mathbf{x}_g \ \mathbf{u}_g]^T$)
- The likelihood of $P(\mathbf{z}|\mathbf{x}, \mathbf{u})$ takes the role of the cost function in typical optimal control
- Conditional distribution $P(\mathbf{u}|\mathbf{x})$ becomes a stochastic control law

$$\max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0:T}, \mathbf{u}_{0:T-1}, \boldsymbol{z}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta}) \prod_{t=0}^{T-1} P(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})$$

- Applies to finite horizon optimal control problems, moving from x_0 to x_g over horizon T
- **z** represents the 'optimization states' and their desired trajectory (i.e. $\mathbf{z}_t = [\mathbf{x}_g \ \mathbf{u}_g]^T$)
- The likelihood of $P(\mathbf{z}|\mathbf{x}, \mathbf{u})$ takes the role of the cost function in typical optimal control
- Conditional distribution $P(\mathbf{u}|\mathbf{x})$ becomes a stochastic control law

$$\max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0:T}, \mathbf{u}_{0:T-1}, \boldsymbol{z}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta}) \prod_{t=0}^{T-1} P(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})$$

- Applies to finite horizon optimal control problems, moving from x_0 to x_g over horizon T
- **z** represents the 'optimization states' and their desired trajectory (i.e. $\mathbf{z}_t = [\mathbf{x}_g \ \mathbf{u}_g]^T$)
- The likelihood of $P(\mathbf{z}|\mathbf{x},\mathbf{u})$ takes the role of the cost function in typical optimal control
- Conditional distribution $P(\mathbf{u}|\mathbf{x})$ becomes a stochastic control law

$$\max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0:T}, \mathbf{u}_{0:T-1}, \boldsymbol{z}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta}) \prod_{t=0}^{T-1} P(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})$$

- Applies to finite horizon optimal control problems, moving from x_0 to x_g over horizon T
- **z** represents the 'optimization states' and their desired trajectory (i.e. $\mathbf{z}_t = [\mathbf{x}_g \ \mathbf{u}_g]^T$)
- The likelihood of $P(\mathbf{z}|\mathbf{x}, \mathbf{u})$ takes the role of the cost function in typical optimal control
- Conditional distribution $P(\mathbf{u}|\mathbf{x})$ becomes a stochastic control law

$$\max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0:T}, \mathbf{u}_{0:T-1}, \boldsymbol{z}_{0:T-1}, \boldsymbol{\theta}) = \max_{\mathbf{u}_{0:T-1}} P(\boldsymbol{x}_{0}) \prod_{t=0}^{T-1} P(\boldsymbol{x}_{t+1} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta}) \prod_{t=0}^{T-1} P(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \mathbf{u}_{t}, \boldsymbol{\theta})$$

- Applies to finite horizon optimal control problems, moving from x_0 to x_g over horizon T
- **z** represents the 'optimization states' and their desired trajectory (i.e. $\mathbf{z}_t = [\mathbf{x}_g \ \mathbf{u}_g]^T$)
- The likelihood of $P(\mathbf{z}|\mathbf{x}, \mathbf{u})$ takes the role of the cost function in typical optimal control
- Conditional distribution $P(\mathbf{u}|\mathbf{x})$ becomes a stochastic control law

LQR:
$$C(\boldsymbol{x}_t, \boldsymbol{u}_t) = (\boldsymbol{x}_g - \boldsymbol{x}_t)^{\mathsf{T}} \mathbf{Q} (\boldsymbol{x}_g - \boldsymbol{x}_t) + (\boldsymbol{u}_g - \boldsymbol{u}_t)^{\mathsf{T}} \mathbf{R} (\boldsymbol{u}_g - \boldsymbol{u}_t)$$
(1)
i2c:
$$\boldsymbol{z}_t = \mathbf{E}_t \boldsymbol{x}_t + \mathbf{E}_t \boldsymbol{u}_t + \mathbf{e}_t + \boldsymbol{\xi}_t$$
(2)

$$\mathbf{z}_t = \mathbf{E}_t \boldsymbol{x}_t + \mathbf{F}_t \mathbf{u}_t + \mathbf{e}_t + \boldsymbol{\xi}_t$$
(2)

$$\boldsymbol{\xi}_t \sim \mathcal{N}(\boldsymbol{0}, \ \boldsymbol{\Sigma}_{\boldsymbol{\xi}}) \tag{3}$$

$$\mathbf{z}_t = \begin{bmatrix} \mathbf{x}_g \\ \mathbf{u}_g \end{bmatrix} \tag{4}$$

$$\boldsymbol{\Sigma}_{\boldsymbol{\xi}} = \frac{1}{\alpha} \begin{bmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix}^{\cdot 1} = \frac{1}{\alpha} \boldsymbol{\Theta}^{\cdot 1}$$
(5)

The Graphical Model



The Graphical Model



Model-based Signal Processing through Message Passing

Linear Gaussian message passing rules are easily constructed from Factor Graph methods⁶:



Where u, Λ denote the 'information' form: $\Lambda = \Sigma^{-1}, \ \nu = \Lambda \mu$

⁶Hans-Andrea Loeliger et al. "The factor graph approach to model-based signal processing". In: Proceedings of the IEEE (2007). Approximate Imput Inference as Stochastic Optimal Control Joe Watson 10/24

As with the LGDS, inference can be performed using EM:

E-Step: Linear Gaussian messsage passing to estimate μ_{x_i}, Σ_{x_i}, μ_{u_i}, Σ_{u_i} from z_{0:T}, μ_{x₀} and Σ_{x₀}. Forward pass: construct priors, i.e μ_{x_i}
Backward pass: improve likelihood, i.e. μ_{x_i}
Marginalisation: calculuate posteriors, i.e. Σ_x =(Λ_{x̄} + Λ_{x̄})⁻¹, μ_x =Σ_x (ν_{x̄} + ν_{x̄})
M-Step: find Σ_ξ, therefore α, via the expected log-likelihood.
update priors on controls μ→ , Σ→

$$\hat{\Sigma}_{\boldsymbol{\xi}} = \frac{1}{T} \sum_{t=0}^{T-1} \left[(\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t}) (\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t})^{\mathsf{T}} + \mathbf{E}_t \boldsymbol{\Sigma}_{x_t} \mathbf{E}_t^{\mathsf{T}} + \mathbf{F}_t \boldsymbol{\Sigma}_{u_t} \mathbf{F}_t^{\mathsf{T}} \right]$$
(6)
$$\alpha = \frac{d_z}{\operatorname{tr}\{\Theta \hat{\Sigma}_{\boldsymbol{\xi}}\}}$$
(7)

$$-\mathcal{L}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}, \alpha) = -\mathcal{L}_{\text{traj}}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}) + \frac{\alpha}{2} \sum_{t=0}^{T-1} C(\boldsymbol{x}_t, \mathbf{u}_t) + \dots$$
(8)

As with the LGDS, inference can be performed using EM:

E-Step: Linear Gaussian messsage passing to estimate μ_{x_i} , Σ_{x_i} , μ_{u_i} , Σ_{u_i} from $\mathbf{z}_{0:T}$, μ_{x_0} and Σ_{x_0} . Forward pass: construct priors, i.e $\mu_{\overrightarrow{x}_i}$ Backward pass: improve likelihood, i.e. $\mu_{\overleftarrow{x}_i}$ Marginalisation: calculuate posteriors, i.e. $\Sigma_x = (\Lambda_{\overrightarrow{x}} + \Lambda_{\overleftarrow{x}})^{-1}$, $\mu_x = \Sigma_x (\nu_{\overrightarrow{x}} + \nu_{\overleftarrow{x}})$ M-Step: find Σ_{ξ} , therefore α , via the expected log-likelihood. update priors on controls $\mu_{\overrightarrow{u}_i}$, $\Sigma_{\overrightarrow{u}_i}$

$$\hat{\Sigma}_{\boldsymbol{\xi}} = \frac{1}{T} \sum_{t=0}^{T-1} \left[(\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t}) (\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t})^{\mathsf{T}} + \mathbf{E}_t \boldsymbol{\Sigma}_{x_t} \mathbf{E}_t^{\mathsf{T}} + \mathbf{F}_t \boldsymbol{\Sigma}_{u_t} \mathbf{F}_t^{\mathsf{T}} \right]$$
(6)
$$\alpha = \frac{d_z}{\operatorname{tr}\{\boldsymbol{\Theta}\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\xi}}\}}$$
(7)

$$-\mathcal{L}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}, \alpha) = -\mathcal{L}_{\text{traj}}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}) + \frac{\alpha}{2} \sum_{t=0}^{T-1} C(\boldsymbol{x}_t, \mathbf{u}_t) + \dots$$
(8)

As with the LGDS, inference can be performed using EM:

- E-Step: Linear Gaussian messsage passing to estimate μ_{x_i} , Σ_{x_i} , μ_{u_i} , Σ_{u_i} from $\mathbf{z}_{0:T}$, μ_{x_0} and Σ_{x_0} . Forward pass: construct priors, i.e $\mu_{\overrightarrow{x}_i}$ Backward pass: improve likelihood, i.e. $\mu_{\overleftarrow{x}_i}$ Marginalisation: calculuate posteriors, i.e. $\Sigma_x = (\Lambda_{\overrightarrow{x}} + \Lambda_{\overleftarrow{x}})^{-1}$, $\mu_x = \Sigma_x (\nu_{\overrightarrow{x}} + \nu_{\overleftarrow{x}})$ M-Step: find Σ_{ε} , therefore α , via the expected log-likelihood.
 - update priors on controls $\mu_{\vec{u}_t}$, $\Sigma_{\vec{u}_t}$

$$\hat{\Sigma}_{\boldsymbol{\xi}} = \frac{1}{T} \sum_{t=0}^{T-1} \left[(\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t}) (\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t})^{\mathsf{T}} + \mathbf{E}_t \boldsymbol{\Sigma}_{x_t} \mathbf{E}_t^{\mathsf{T}} + \mathbf{F}_t \boldsymbol{\Sigma}_{u_t} \mathbf{F}_t^{\mathsf{T}} \right]$$
(6)
$$\alpha = \frac{d_z}{\operatorname{tr}\{\boldsymbol{\Theta} \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\xi}}\}}$$
(7)

$$-\mathcal{L}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}, \alpha) = -\mathcal{L}_{\text{traj}}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}) + \frac{\alpha}{2} \sum_{t=0}^{T-1} C(\boldsymbol{x}_t, \mathbf{u}_t) + \dots$$
(8)

As with the LGDS, inference can be performed using EM:

- E-Step: Linear Gaussian messsage passing to estimate μ_{x_i} , Σ_{x_i} , μ_{u_i} , Σ_{u_i} from $\mathbf{z}_{0:T}$, μ_{x_0} and Σ_{x_0} . Forward pass: construct priors, i.e $\mu_{\overrightarrow{x}_i}$ Backward pass: improve likelihood, i.e. $\mu_{\overleftarrow{x}_i}$ Marginalisation: calculuate posteriors, i.e. $\Sigma_x = (\Lambda_{\overrightarrow{x}} + \Lambda_{\overleftarrow{x}})^{-1}$, $\mu_x = \Sigma_x (\nu_{\overrightarrow{x}} + \nu_{\overleftarrow{x}})$ M-Step: find Σ_{ε} , therefore α , via the expected log-likelihood.
- where α_{ξ} , therefore α_{ξ} , via the expected log-likelihood. update priors on controls $\mu_{\vec{u}_{t}}$, $\Sigma_{\vec{u}_{t}}$

$$\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\xi}} = \frac{1}{T} \sum_{t=0}^{T-1} \left[(\mathbf{z}_t - \mathbf{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t}) (\mathbf{z}_t - \boldsymbol{E}_t \boldsymbol{\mu}_{x_t} - \mathbf{F}_t \boldsymbol{\mu}_{u_t})^{\mathsf{T}} + \boldsymbol{E}_t \boldsymbol{\Sigma}_{x_t} \mathbf{E}_t^{\mathsf{T}} + \mathbf{F}_t \boldsymbol{\Sigma}_{u_t} \mathbf{F}_t^{\mathsf{T}} \right]$$
(6)
$$\alpha = \frac{d_z}{\operatorname{tr}\{\boldsymbol{\Theta}\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\xi}}\}}$$
(7)

$$-\mathcal{L}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}, \alpha) = -\mathcal{L}_{\text{traj}}(\boldsymbol{x}_{0:T-1}, \mathbf{u}_{0:T-1}) + \frac{\alpha}{2} \sum_{t=0}^{T-1} C(\boldsymbol{x}_t, \mathbf{u}_t) + \dots$$
(8)

Experiment: Simple Linear System



Figure 4: Demonstrating how i2c generalizes the Dynamic Programming Finite Horizon LQR trajectory. Note 'Filtered' and 'Prediction' correspond to $\mu_{\vec{x}'_{i+1}}$ and $\mu_{\vec{x}'_{i+1}}$.

$$P_{t} = Q + A^{T} P_{t+1} A - A^{T} P_{t+1} B (R + B^{T} P_{t+1} B)^{-1} B^{T} P_{t+1} A$$
(9)

$$\boldsymbol{\Lambda}_{\boldsymbol{\widehat{x}}_{t}} = \boldsymbol{E}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \boldsymbol{F}_{t} \boldsymbol{\Sigma}_{\boldsymbol{\overrightarrow{u}}_{t}} \boldsymbol{F}_{t}^{\mathsf{T}})^{-1} \boldsymbol{E}_{t} + \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\boldsymbol{\widehat{x}}_{t+1}} \boldsymbol{A}_{t} \\
-\boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\boldsymbol{\overrightarrow{x}}_{t+1}} ((\boldsymbol{\Sigma}_{\eta_{t}} + \boldsymbol{B}_{t} (\boldsymbol{\Lambda}_{\boldsymbol{\overrightarrow{u}}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \boldsymbol{E}_{t} \boldsymbol{\Sigma}_{\boldsymbol{\overrightarrow{x}}_{t}} \boldsymbol{E}_{t}^{\mathsf{T}})^{-1} \boldsymbol{F}_{t})^{-1} \boldsymbol{B}_{t}^{\mathsf{T}})^{-1} + \boldsymbol{\Lambda}_{\boldsymbol{\overleftarrow{x}}_{t+1}})^{-1} \boldsymbol{\Lambda}_{\boldsymbol{\overleftarrow{x}}_{t+1}} \boldsymbol{A}_{t} \quad (10)$$

$$P_{t} = Q + A^{T} P_{t+1} A - A^{T} P_{t+1} B (R + B^{T} P_{t+1} B)^{-1} B^{T} P_{t+1} A$$
(9)

$$\boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t}} = \boldsymbol{E}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{F}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{u}}_{t}} \mathbf{F}_{t}^{\mathsf{T}})^{-1} \boldsymbol{E}_{t} + \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}} \boldsymbol{A}_{t} - \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}} ((\boldsymbol{\Sigma}_{\boldsymbol{\eta}_{t}} + \boldsymbol{B}_{t} (\boldsymbol{\Lambda}_{\overrightarrow{\boldsymbol{u}}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{x}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{-1} \mathbf{B}_{t}^{\mathsf{T}})^{-1} + \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}})^{-1} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}} \boldsymbol{A}_{t}$$
(10)

$$P_{t} = Q + A^{T} P_{t+1} A - A^{T} P_{t+1} B (R + B^{T} P_{t+1} B)^{-1} B^{T} P_{t+1} A$$
(9)

$$\boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t}} = \boldsymbol{E}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{F}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{u}}_{t}} \mathbf{F}_{t}^{\mathsf{T}})^{-1} \boldsymbol{E}_{t} + \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}} \boldsymbol{A}_{t}
- \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}} ((\boldsymbol{\Sigma}_{\boldsymbol{\eta}_{t}} + \boldsymbol{B}_{t} (\boldsymbol{\Lambda}_{\overrightarrow{\boldsymbol{u}}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{x}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{-1} \mathbf{B}_{t}^{\mathsf{T}})^{-1} + \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}})^{-1} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{x}}_{t+1}} \boldsymbol{A}_{t}$$
(10)

$$P_{t} = Q + A^{T} P_{t+1} A - A^{T} P_{t+1} B (R + B^{T} P_{t+1} B)^{-1} B^{T} P_{t+1} A$$
(9)

$$\boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t}} = \boldsymbol{E}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{F}_{t} \boldsymbol{\Sigma}_{\overrightarrow{u}_{t}} \mathbf{F}_{t}^{\mathsf{T}})^{-1} \boldsymbol{E}_{t} + \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{A}_{t} - \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} ((\boldsymbol{\Sigma}_{\boldsymbol{\eta}_{t}} + \boldsymbol{B}_{t} (\boldsymbol{\Lambda}_{\overrightarrow{u}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\mathbf{x}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{-1} \mathbf{B}_{t}^{\mathsf{T}})^{-1} + \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}})^{-1} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{A}_{t}$$
(10)

$$P_{t} = Q + A^{T} P_{t+1} A - A^{T} P_{t+1} B (R + B^{T} P_{t+1} B)^{-1} B^{T} P_{t+1} A$$
(9)

$$\boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t}} = \boldsymbol{E}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{F}_{t} \boldsymbol{\Sigma}_{\overrightarrow{u}_{t}} \mathbf{F}_{t}^{\mathsf{T}})^{-1} \boldsymbol{E}_{t} + \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{A}_{t} - \boldsymbol{A}_{t}^{\mathsf{T}} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} ((\boldsymbol{\Sigma}_{\boldsymbol{\eta}_{t}} + \boldsymbol{B}_{t} (\boldsymbol{\Lambda}_{\overrightarrow{u}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\mathbf{x}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{-1} \boldsymbol{B}_{t}^{\mathsf{T}})^{-1} + \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}})^{-1} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{A}_{t}$$
(10)

Like LQR, can i2c give us time-varying linear controllers?

The E-Step gave us the joint Normal distributions $P(\boldsymbol{x}_t, \mathbf{u}_t)$. The conditional

$$P(\mathbf{u}_t | \boldsymbol{x}_t) : \boldsymbol{\mu}_{u_t | \boldsymbol{x}_t} = \boldsymbol{\mu}_{u_t} + \boldsymbol{\Sigma}_{u_t \boldsymbol{x}_t} \boldsymbol{\Sigma}_{\boldsymbol{x}_t \boldsymbol{x}_t}^{-1} (\boldsymbol{x}_t - \boldsymbol{\mu}_{\boldsymbol{x}_t}),$$
(11)

has a mean linear is x.

Therefore we get time-varying linear Gaussian controllers

$$P(\mathbf{u}_t | \boldsymbol{x}_t) = \mathcal{N}(\mathbf{K}_t \boldsymbol{x}_t + \mathbf{k}_t, \boldsymbol{\Sigma}_{k_t})$$
(12)

Maximum Entropy LQR:

$$\boldsymbol{K}_{t} = -(\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1} \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{A}$$
(13)

$$\mathbf{k}_{t} = (\mathbf{R} + \mathbf{B}^{T} \mathbf{P}_{t+1} \mathbf{B})^{-1} (\mathbf{B}^{T} (\mathbf{P}_{t+1} \mathbf{a}_{t} + \mathbf{p}_{t+1}) - \mathbf{R} \mathbf{u}_{g})$$
(14)

$$\boldsymbol{\Sigma}_{k_t} = (\boldsymbol{R} + \boldsymbol{B}^T \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1}$$
⁽¹⁵⁾

Expressing $P(\mathbf{u}|\boldsymbol{x})$ through the messages:

$$K_{t} = -\Sigma_{u_{t}} B_{t} \Gamma_{t+1} A_{\overleftarrow{\Sigma}_{t+1}} \Psi_{t+1} A_{t}, \qquad (16)$$

$$k_{t} = \Sigma_{u_{t}} (\nu_{\overrightarrow{u}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\Sigma_{\boldsymbol{\xi}} \mathbf{E}_{t} \Sigma_{\overrightarrow{\boldsymbol{\chi}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{-1} (\mathbf{z}_{t} - \mathbf{E}_{t} \mu_{\overrightarrow{\boldsymbol{\chi}}_{t}} - \mathbf{e}_{t}) + B_{t}^{\mathsf{T}} (\Gamma_{t+1} \nu_{\overleftarrow{\Sigma}_{t+1}} + (\mathbf{I} - \Gamma_{t+1}) \nu_{\overrightarrow{\boldsymbol{\chi}}_{t}''} - \Gamma_{t+1} A_{\overleftarrow{\Sigma}_{t+1}} \Psi_{t+1} \mathbf{a}_{t})), \qquad (17)$$

$$\Sigma_{k_{t}} = \Sigma_{u_{t}} = (A_{\overrightarrow{\boldsymbol{u}}} + \mathbf{F}_{t}^{\mathsf{T}} (\Sigma_{\boldsymbol{\xi}} + \mathbf{E}_{t} \Sigma_{\overrightarrow{\boldsymbol{\chi}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{-1} \mathbf{F}_{t} + B_{t}^{\mathsf{T}} \Gamma_{t+1} A_{\overleftarrow{\Sigma}_{t+1}} B_{t})^{-1} \qquad (18)$$

Maximum Entropy LQR:

$$\boldsymbol{K}_{t} = -(\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1} \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{A}$$
(13)

$$\mathbf{k}_{t} = (\mathbf{R} + \mathbf{B}^{T} \mathbf{P}_{t+1} \mathbf{B})^{-1} (\mathbf{B}^{T} (\mathbf{P}_{t+1} \mathbf{a}_{t} + \mathbf{p}_{t+1}) - \mathbf{R} \mathbf{u}_{g})$$
(14)

$$\boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} = (\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1}$$
(15)

$$\begin{aligned} \mathbf{K}_{t} &= -\boldsymbol{\Sigma}_{u_{t}} \boldsymbol{B}_{t} \boldsymbol{\Gamma}_{t+1} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{\Psi}_{t+1} \boldsymbol{A}_{t}, \end{aligned} \tag{16} \\ \mathbf{k}_{t} &= \boldsymbol{\Sigma}_{u_{t}} \left(\boldsymbol{\nu}_{\overrightarrow{u}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} \left(\boldsymbol{\Sigma}_{\boldsymbol{\xi}} \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\mathbf{x}}_{t}} \mathbf{E}_{t}^{\mathsf{T}} \right)^{-1} \left(\mathbf{z}_{t} - \mathbf{E}_{t} \boldsymbol{\mu}_{\overrightarrow{\mathbf{x}}_{t}} - \mathbf{e}_{t} \right) \\ &+ \boldsymbol{B}_{t}^{\mathsf{T}} \left(\boldsymbol{\Gamma}_{t+1} \boldsymbol{\nu}_{\overleftarrow{\mathbf{x}}_{t+1}} + (\mathbf{I} - \boldsymbol{\Gamma}_{t+1}) \boldsymbol{\nu}_{\overrightarrow{\mathbf{x}}_{t}''} - \boldsymbol{\Gamma}_{t+1} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{\Psi}_{t+1} \mathbf{a}_{t} \right) \right), \end{aligned} \tag{17} \\ \boldsymbol{\Sigma}_{k_{t}} &= \boldsymbol{\Sigma}_{u_{t}} = \left(\boldsymbol{\Lambda}_{\overrightarrow{u}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} \left(\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\mathbf{x}}_{t}} \mathbf{E}_{t}^{\mathsf{T}} \right)^{-1} \mathbf{F}_{t} + \boldsymbol{B}_{t}^{\mathsf{T}} \boldsymbol{\Gamma}_{t+1} \boldsymbol{\Lambda}_{\overleftarrow{\mathbf{x}}_{t+1}} \boldsymbol{B}_{t} \right)^{-1} \end{aligned} \tag{18}$$

Maximum Entropy LQR:

$$\boldsymbol{K}_{t} = -(\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1} \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{A}$$
(13)

$$\mathbf{k}_{t} = (\mathbf{R} + \mathbf{B}^{T} \mathbf{P}_{t+1} \mathbf{B})^{-1} (\mathbf{B}^{T} (\mathbf{P}_{t+1} \mathbf{a}_{t} + \mathbf{p}_{t+1}) - \mathbf{R} \mathbf{u}_{g})$$
(14)

$$\boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} = (\boldsymbol{R} + \boldsymbol{B}^{\mathrm{T}} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1}$$
(15)

$$\begin{aligned} \boldsymbol{K}_{t} &= -\boldsymbol{\Sigma}_{\boldsymbol{u}_{t}}\boldsymbol{B}_{t}\boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{\chi}}_{t+1}}\boldsymbol{\Psi}_{t+1}\boldsymbol{A}_{t}, \end{aligned} \tag{16} \\ \boldsymbol{k}_{t} &= \boldsymbol{\Sigma}_{\boldsymbol{u}_{t}}(\boldsymbol{\nu}_{\overrightarrow{\boldsymbol{u}}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}}(\boldsymbol{\Sigma}_{\boldsymbol{\xi}}\boldsymbol{E}_{t}\boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{\chi}}_{t}}\boldsymbol{E}_{t}^{\mathsf{T}})^{-1}(\boldsymbol{z}_{t} - \boldsymbol{E}_{t}\boldsymbol{\mu}_{\overrightarrow{\boldsymbol{\chi}}_{t}} - \boldsymbol{e}_{t}) \\ &+ \boldsymbol{B}_{t}^{\mathsf{T}}(\boldsymbol{\Gamma}_{t+1}\boldsymbol{\nu}_{\overleftarrow{\boldsymbol{\chi}}_{t+1}} + (\boldsymbol{I} - \boldsymbol{\Gamma}_{t+1})\boldsymbol{\nu}_{\overrightarrow{\boldsymbol{\chi}}_{t}''} - \boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{\chi}}_{t+1}}\boldsymbol{\Psi}_{t+1}\boldsymbol{a}_{t})), \end{aligned} \tag{17} \\ \boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} &= \boldsymbol{\Sigma}_{\boldsymbol{u}_{t}} = (\boldsymbol{\Lambda}_{\overrightarrow{\boldsymbol{u}}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}}(\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \boldsymbol{E}_{t}\boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{\chi}}_{t}}\boldsymbol{E}_{t}^{\mathsf{T}})^{-1}\boldsymbol{F}_{t} + \boldsymbol{B}_{t}^{\mathsf{T}}\boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{\chi}}_{t+1}}\boldsymbol{B}_{t})^{-1} \end{aligned} \tag{18}$$

Maximum Entropy LQR:

$$\boldsymbol{K}_{t} = -(\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1} \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{A}$$
(13)

$$\mathbf{k}_{t} = (\mathbf{R} + \mathbf{B}^{T} \mathbf{P}_{t+1} \mathbf{B})^{-1} (\mathbf{B}^{T} (\mathbf{P}_{t+1} \mathbf{a}_{t} + \mathbf{p}_{t+1}) - \mathbf{R} \mathbf{u}_{g})$$
(14)

$$\boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} = (\boldsymbol{R} + \boldsymbol{B}^{\mathrm{T}} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1}$$
(15)

$$\begin{split} \boldsymbol{K}_{t} &= -\boldsymbol{\Sigma}_{\boldsymbol{u}_{t}} \boldsymbol{B}_{t} \boldsymbol{\Gamma}_{t+1} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{\Sigma}}_{t+1}} \boldsymbol{\Psi}_{t+1} \boldsymbol{A}_{t}, \end{split}$$
(16)
$$\boldsymbol{k}_{t} &= \boldsymbol{\Sigma}_{\boldsymbol{u}_{t}} (\boldsymbol{\nu}_{\overrightarrow{\boldsymbol{u}}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{\chi}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{\cdot 1} (\mathbf{z}_{t} - \mathbf{E}_{t} \boldsymbol{\mu}_{\overrightarrow{\boldsymbol{\chi}}_{t}} - \mathbf{e}_{t}) \\ &+ \boldsymbol{B}_{t}^{\mathsf{T}} (\boldsymbol{\Gamma}_{t+1} \boldsymbol{\nu}_{\overleftarrow{\boldsymbol{\Sigma}}_{t+1}} + (\mathbf{I} - \boldsymbol{\Gamma}_{t+1}) \boldsymbol{\nu}_{\overrightarrow{\boldsymbol{\chi}}_{t}''} - \boldsymbol{\Gamma}_{t+1} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{\Sigma}}_{t+1}} \boldsymbol{\Psi}_{t+1} \mathbf{a}_{t})), \end{aligned}$$
(17)
$$\boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} &= \boldsymbol{\Sigma}_{\boldsymbol{u}_{t}} = (\boldsymbol{\Lambda}_{\overrightarrow{\boldsymbol{u}}_{t}} + \mathbf{F}_{t}^{\mathsf{T}} (\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \mathbf{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\boldsymbol{\chi}}_{t}} \mathbf{E}_{t}^{\mathsf{T}})^{\cdot 1} \mathbf{F}_{t} + \boldsymbol{B}_{t}^{\mathsf{T}} \boldsymbol{\Gamma}_{t+1} \boldsymbol{\Lambda}_{\overleftarrow{\boldsymbol{\chi}}_{t+1}} \boldsymbol{B}_{t})^{-1} \end{aligned}$$
(18)

Maximum Entropy LQR:

$$\boldsymbol{K}_{t} = -(\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1} \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{A}$$
(13)

$$\mathbf{k}_{t} = (\mathbf{R} + \mathbf{B}^{T} \mathbf{P}_{t+1} \mathbf{B})^{-1} (\mathbf{B}^{T} (\mathbf{P}_{t+1} \mathbf{a}_{t} + \mathbf{p}_{t+1}) - \mathbf{R} \mathbf{u}_{g})$$
(14)

$$\boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} = (\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1}$$
(15)

$$\begin{aligned} \boldsymbol{K}_{t} &= -\boldsymbol{\Sigma}_{u_{t}}\boldsymbol{B}_{t}\boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\boldsymbol{\widehat{\chi}}_{t+1}}\boldsymbol{\Psi}_{t+1}\boldsymbol{A}_{t}, \end{aligned} \tag{16} \\ \boldsymbol{k}_{t} &= \boldsymbol{\Sigma}_{u_{t}}(\boldsymbol{\nu}_{\overrightarrow{u}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}}(\boldsymbol{\Sigma}_{\boldsymbol{\xi}}\boldsymbol{E}_{t}\boldsymbol{\Sigma}_{\boldsymbol{\widehat{\chi}}_{t}}\boldsymbol{E}_{t}^{\mathsf{T}})^{\cdot1}(\boldsymbol{z}_{t} - \boldsymbol{E}_{t}\boldsymbol{\mu}_{\boldsymbol{\widehat{\chi}}_{t}} - \boldsymbol{e}_{t}) \\ &+ \boldsymbol{B}_{t}^{\mathsf{T}}(\boldsymbol{\Gamma}_{t+1}\boldsymbol{\nu}_{\boldsymbol{\widehat{\chi}}_{t+1}} + (\boldsymbol{I} - \boldsymbol{\Gamma}_{t+1})\boldsymbol{\nu}_{\boldsymbol{\widehat{\chi}}_{t}''} - \boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\boldsymbol{\widehat{\chi}}_{t+1}}\boldsymbol{\Psi}_{t+1}\boldsymbol{a}_{t})), \end{aligned} \tag{17} \\ \boldsymbol{\Sigma}_{k_{t}} &= \boldsymbol{\Sigma}_{u_{t}} = (\boldsymbol{\Lambda}_{\overrightarrow{u}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}}(\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \boldsymbol{E}_{t}\boldsymbol{\Sigma}_{\boldsymbol{\widehat{\chi}}_{t}}\boldsymbol{E}_{t}^{\mathsf{T}})^{\cdot1}\boldsymbol{F}_{t} + \boldsymbol{B}_{t}^{\mathsf{T}}\boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\boldsymbol{\widehat{\chi}}_{t+1}}\boldsymbol{B}_{t})^{-1} \end{aligned} \tag{18}$$

Maximum Entropy LQR:

$$\boldsymbol{K}_{t} = -(\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1} \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{A}$$
(13)

$$\mathbf{k}_{t} = (\mathbf{R} + \mathbf{B}^{T} \mathbf{P}_{t+1} \mathbf{B})^{-1} (\mathbf{B}^{T} (\mathbf{P}_{t+1} \mathbf{a}_{t} + \mathbf{p}_{t+1}) - \mathbf{R} \mathbf{u}_{g})$$
(14)

$$\boldsymbol{\Sigma}_{\boldsymbol{k}_{t}} = (\boldsymbol{R} + \boldsymbol{B}^{T} \boldsymbol{P}_{t+1} \boldsymbol{B})^{-1}$$
(15)

$$\begin{aligned} \boldsymbol{K}_{t} &= -\boldsymbol{\Sigma}_{u_{t}}\boldsymbol{B}_{t}\boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\boldsymbol{\widehat{\chi}}_{t+1}}\boldsymbol{\Psi}_{t+1}\boldsymbol{A}_{t}, \end{aligned} \tag{16} \\ \boldsymbol{k}_{t} &= \boldsymbol{\Sigma}_{u_{t}} \left(\boldsymbol{\nu}_{\overrightarrow{u}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}} \left(\boldsymbol{\Sigma}_{\boldsymbol{\xi}} \boldsymbol{E}_{t} \boldsymbol{\Sigma}_{\overrightarrow{\chi}_{t}} \boldsymbol{E}_{t}^{\mathsf{T}}\right)^{\cdot 1} (\boldsymbol{z}_{t} - \boldsymbol{E}_{t}\boldsymbol{\mu}_{\overrightarrow{\chi}_{t}} - \boldsymbol{e}_{t}) \\ &+ \boldsymbol{B}_{t}^{\mathsf{T}} \left(\boldsymbol{\Gamma}_{t+1}\boldsymbol{\nu}_{\overrightarrow{\chi}_{t+1}} + (\boldsymbol{I} - \boldsymbol{\Gamma}_{t+1})\boldsymbol{\nu}_{\overrightarrow{\chi}_{t}''} - \boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\overleftarrow{\chi}_{t+1}}\boldsymbol{\Psi}_{t+1}\boldsymbol{a}_{t}\right)), \end{aligned} \tag{17} \\ \boldsymbol{\Sigma}_{k_{t}} &= \boldsymbol{\Sigma}_{u_{t}} = \left(\boldsymbol{\Lambda}_{\overrightarrow{u}_{t}} + \boldsymbol{F}_{t}^{\mathsf{T}} \left(\boldsymbol{\Sigma}_{\boldsymbol{\xi}} + \boldsymbol{E}_{t}\boldsymbol{\Sigma}_{\overrightarrow{\chi}_{t}} \boldsymbol{E}_{t}^{\mathsf{T}}\right)^{\cdot 1} \boldsymbol{F}_{t} + \boldsymbol{B}_{t}^{\mathsf{T}} \boldsymbol{\Gamma}_{t+1}\boldsymbol{\Lambda}_{\overleftarrow{\chi}_{t+1}} \boldsymbol{B}_{t}\right)^{-1} \end{aligned} \tag{18}$$

Experiment: Simple Linear System



Figure 5: Demonstrating how i2c generalizes the Dynamic Programming Finite Horizon LQR controller.

i2c for Nonlinear Systems

Taking inspiration from the Extended Kalman Filter, we **linearize the dynamics** about the state action trajectory during the forward pass.

However, this local linearization assumption reduces the amount of uncertainty in the trajectory we can tolerate.

As the expected log-likelihood in the M-step is now approximate, we must regularize of the $\Sigma_{\xi}(\alpha)$ update. Starting with a KL bound, this can be practically implemented as a bound of the update ratio:

$$\frac{\alpha^i}{\alpha^{i+1}} \le \delta_\alpha$$





The Algorithm

Data: Task: {T, \mathbf{z}_{g} , Θ , $\mu_{\overrightarrow{x}_{\alpha}}$, $\mathbf{f}(\mathbf{x}, \mathbf{u})$, $g(\mathbf{x}, \mathbf{u})$ }, Priors: { $\Sigma_{\overrightarrow{x}_{\alpha}}$, $\mu_{\overrightarrow{u}_{\alpha}}$, $\Sigma_{\overrightarrow{u}_{\alpha}}$ }, Hyperparameters: { α, δ_{α} } **Result:** \mathbf{K}_t , \mathbf{k}_t , for t = 0 : T-1while not converged do // E-Step for $i \leftarrow 0$ to T - 1 do Compute $\mu_{\vec{x}_{t+1}}$, $\Sigma_{\vec{x}_{t+1}}$ from forward messages, updating A_t , a_t , B_t , E_t , e_t and F_t end for $i \leftarrow T$ to 1 do Compute μ_x , Σ_x , μ_y , Σ_y , from backward messages and marginalisation end // M-Step Update α with regularization, update priors $\mu_{\vec{u}_{\alpha,r}} = \mu_{u_{\alpha,r}}, \Sigma_{\vec{u}_{\alpha,r}} = \Sigma_{u_{\alpha,r}}$ end 11 Controller Computer linear Gaussian controller $\mathbf{K}_t, \mathbf{k}_t, \mathbf{\Sigma}_{k_t}$ for t = 0:T-1 from messages

Algorithm 1: EM for Linearized Gaussian i2c



We compared i2c against Iterative LQR (iLQR)⁷ and Guided Policy Search (GPS)⁸



Figure 6: Comparison of the trajectory cost prediction over iterations during trajectory optimization.

⁷Weiwei Li and Emanuel Todorov. "Iterative linear quadratic regulator design for nonlinear biological movement systems.". In: *ICINCO* (1). 2004.

⁸Sergey Levine and Vladlen Koltun. "Guided Policy Search". In: Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28. ICML'13. 2013.

Environment	Algorithm	Predicted Cost	Evaluated Cost
Pendulum	i2c iLQR GPS	$egin{array}{llllllllllllllllllllllllllllllllllll$	$\begin{array}{c} {\bf 1.37 \times 10^4 \pm 3.82} \\ {1.11 \times 10^5 \pm 20.38} \\ {7.01 \times 10^4 \pm 30.96} \end{array}$
Cartpole	i2c iLQR GPS	1.73×10^5 1.14×10^5 1.65×10^5	$\begin{array}{c} {\bf 1.74 \times 10^5 \pm 0.14} \\ {\bf 1.76 \times 10^7 \pm 88.63} \\ {\bf 2.94 \times 10^6 \pm 17.60} \end{array}$
Double Cartpole	i2c iLQR GPS	3.12×10^5 2.37×10^5 3.76×10^5	$\begin{array}{c} {\bf 3.21 \times 10^5 \pm 1.79} \\ {1.76 \times 10^7 \pm 5.27 \times 10^5} \\ {2.94 \times 10^6 \pm 44.39} \end{array}$

Table 1: Evaluating the optimized deterministic controller of each algorithm on the simulated stochastic environments. Predicted Cost refers to the converged value from Figure 6, Evaluated Cost shows the mean and standard deviation after 100 trials.



Animations

i2c Pendulum

iLQR Pendulum

Conclusion

• LQR can be framed as input inference of a linear Gaussian dynamical system

- i2c represents a practical **EM algorithm** for performing Bayesian nonlinear trajectory optimization through **message passing** and approximate inference, returning time-varying linear controllers
- Through use of Bayesian inference, the scheme naturally exhibits **uncertainty-derived regularization** that has been explicitly incorporated into (or lacking in) previous methods

Conclusion

- LQR can be framed as input inference of a linear Gaussian dynamical system
- i2c represents a practical EM algorithm for performing Bayesian nonlinear trajectory optimization through **message passing** and approximate inference, returning time-varying linear controllers
- Through use of Bayesian inference, the scheme naturally exhibits **uncertainty-derived regularization** that has been explicitly incorporated into (or lacking in) previous methods

Conclusion

- LQR can be framed as input inference of a linear Gaussian dynamical system
- i2c represents a practical EM algorithm for performing Bayesian nonlinear trajectory optimization through **message passing** and approximate inference, returning time-varying linear controllers
- Through use of Bayesian inference, the scheme naturally exhibits **uncertainty-derived regularization** that has been explicitly incorporated into (or lacking in) previous methods

Questions?

- Boedecker, Joschka et al. "Approximate real-time optimal control based on sparse gaussian process models". In: 2014 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL). IEEE. 2014.
- Chua, Kurtland et al. "Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models". In: Advances in Neural Information Processing Systems 31.
- Deisenroth, Marc and Carl E Rasmussen. "PILCO: A model-based and data-efficient approach to policy search". In: Proceedings of the 28th International Conference on machine learning (ICML-11). 2011.
- Hennig, Philipp, Michael A Osborne, and Mark Girolami. "Probabilistic numerics and uncertainty in computations". In: Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences (2015).
- Hoffmann, Christian and Philipp Rostalski. "Linear Optimal Control on Factor Graphs A Message Passing Perspective -". In: IFAC (International Federation of Automatic Control) (2017).
- Levine, Sergey. "Motor skill learning with local trajectory methods". PhD thesis. Stanford University, 2014.
- Levine, Sergey and Vladlen Koltun. "Guided Policy Search". In: Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28. ICML'13. 2013.
- Li, Weiwei and Emanuel Todorov. "Iterative linear quadratic regulator design for nonlinear biological movement systems.". In: *ICINCO* (1). 2004.
- Loeliger, Hans-Andrea et al. "The factor graph approach to model-based signal processing". In: *Proceedings of the IEEE* (2007).
 - Toussaint, Marc. "Robot trajectory optimization using approximate inference". In: *Proceedings of the 26th annual international conference on machine learning*. ACM. 2009.



Previous Work

	Controller?	Adaptive Input Uncertainty?	Nonlinear?	Uncalibrated Cost?
AICO ⁹	X	×	\checkmark	×
Hoffmann <i>et al</i> ¹⁰	\checkmark	×	X	X
i2c	\checkmark	\checkmark	\checkmark	\checkmark

⁹Marc Toussaint. "Robot trajectory optimization using approximate inference". In: Proceedings of the 26th annual international conference on machine learning. ACM. 2009.

¹⁰Christian Hoffmann and Philipp Rostalski. "Linear Optimal Control on Factor Graphs - A Message Passing Perspective -". In: IFAC (International Federation of Automatic Control) (2017).



Fig. 3. Factor graph for LQG control.