

# Model-Based and Model-Free Imitation Learning

---

**Hakan Girgin**

25.09.2020



# Learning from demonstration (LfD) in robotics

$\mathbf{x}_t$ : **state** (e.g. position and velocity) of the robot at timestep  $t$

$\mathbf{u}_t$ : **control command** applied by the robot at timestep  $t$

$\mathcal{S}$ : **context variable(s)** such as initial/final positions, user preferences and environment properties.

**Dataset:**  $\mathcal{D} = \{\mathbf{x}_t, \mathbf{u}_t\}_{t=0}^T, \mathcal{S}$

Trajectory level abstraction

$$p(\boldsymbol{\tau} | \mathbf{x}_0, \mathcal{S})$$

State level abstraction

$$p(\mathbf{x}_t | t, \mathcal{S}), \quad p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathcal{S})$$

Action-State abstraction

$$p(\mathbf{u}_t | \mathbf{x}_t, \mathcal{S})$$

DMP

GMM-GMR

HMM-LQR

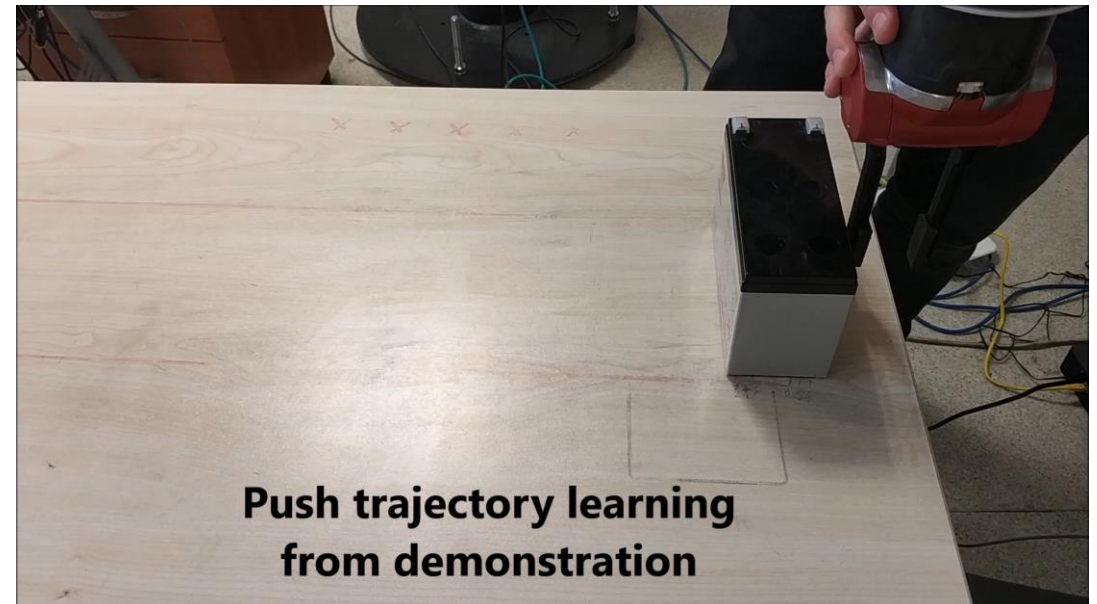
SEDS

ProMP

End-to-end deep  
visuomotor

KMP

BGMM



# Learning from demonstration (LfD) in robotics

$\mathbf{x}_t$ : **state** (e.g. position and velocity) of the robot at timestep  $t$

$\mathbf{u}_t$ : **control command** applied by the robot at timestep  $t$

$\mathcal{S}$ : **context variable(s)** such as initial/final positions, user preferences and environment properties.

**Dataset:**  $\mathcal{D} = \{\mathbf{x}_t, \mathbf{u}_t\}_{t=0}^T, \mathcal{S}$

## Trajectory level abstraction

$$p(\boldsymbol{\tau} | \mathbf{x}_0, \mathcal{S})$$

## State level abstraction

$$p(\mathbf{x}_t | t, \mathcal{S}), \quad p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathcal{S})$$

## Action-State abstraction

$$p(\mathbf{u}_t | \mathbf{x}_t, \mathcal{S})$$

DMP

GMM-GMR

HMM-LQR

SEDS

ProMP

End-to-end deep  
visuomotor

KMP

BGMM

Model-free  
LfD

**Note:** They can still be combined with forward-models in the control level. But the learning is not model-aware.

# Model-aware learning $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$

What are the problems that might arise if we do not care about the forward-model of the robot in learning movement primitives?

- ...
- Generalization to infeasible trajectories
- If not kinesthetic teaching, correspondence problems
- Causality
- High gains in the controller



Model-aware learning:

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t$$

$$p(\mathbf{x}_{t+1}) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t$$

$$p(\boldsymbol{\tau}) = p(\mathbf{x}_0) \prod_{t=1}^T p(\mathbf{x}_{t+1}|\mathbf{x}_t)$$

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t)$$

$$p(\mathbf{x}_t|t)$$

$$p(\boldsymbol{\tau}|\mathbf{x}_0)$$

$$p(\mathbf{u}_t|\mathbf{x}_t)$$

# Probabilistic model-based imitation learning\*

Matching :

- Probability density of the trajectory calculated from demonstration.  $p_{\text{exp}}(\boldsymbol{\tau})$
- Probability density of the trajectory parametrized by the unknown policy parameters  $p_{\boldsymbol{\theta}}(\boldsymbol{\tau})$

$$\text{KL}(p(x)||q(x)) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

$$\min_{\boldsymbol{\theta}} \text{KL}(p_{\text{exp}}(\boldsymbol{\tau})||p_{\boldsymbol{\theta}}(\boldsymbol{\tau}))$$

➔  $\min_{\boldsymbol{\theta}} \sum_{t=1}^T \text{KL}(p_{\text{exp}}(\mathbf{x}_t)||p_{\boldsymbol{\theta}}(\mathbf{x}_t))$

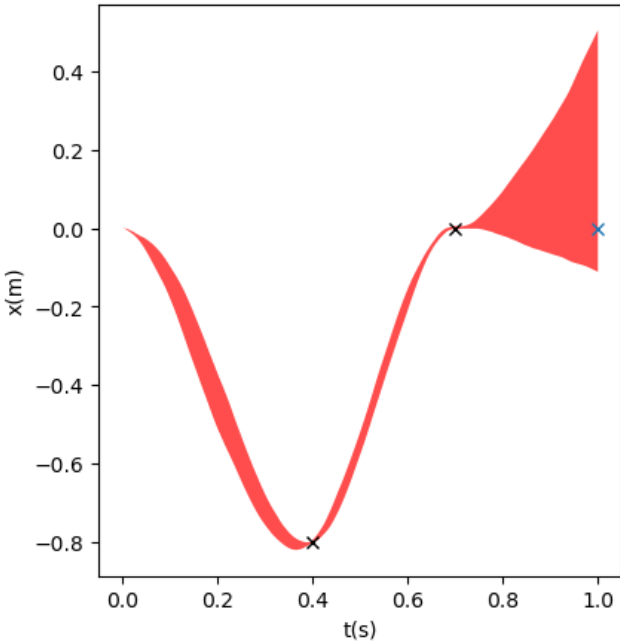
Assumption 1

$$p(\boldsymbol{\tau}) = \prod_{t=1}^T p(\mathbf{x}_t)$$

# Probabilistic model-based imitation learning\*

How to calculate  $p_{\text{exp}}(\boldsymbol{\tau})$  ?

Demonstrations



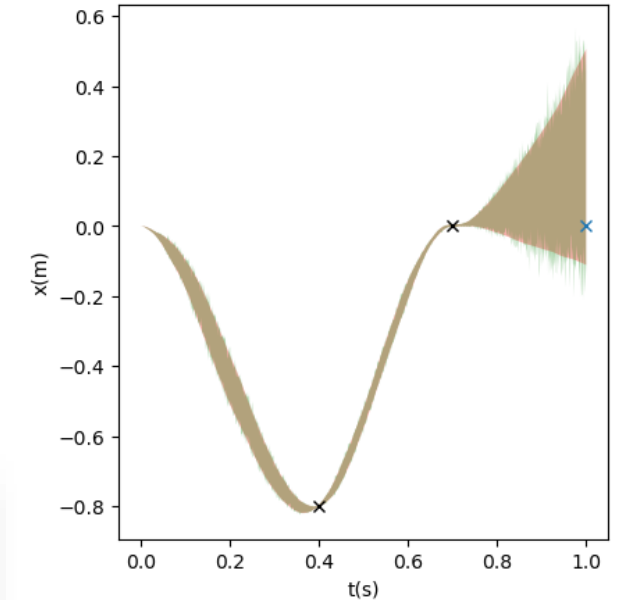
$$\hat{\boldsymbol{\mu}}_t^{\text{exp}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_t^i$$

$$\hat{\boldsymbol{\Sigma}}_t^{\text{exp}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_t^i - \hat{\boldsymbol{\mu}}_t^{\text{exp}})(\mathbf{x}_t^i - \hat{\boldsymbol{\mu}}_t^{\text{exp}})^\top$$

$$p(\boldsymbol{\tau}^{\text{exp}}) = \mathcal{N}(\hat{\boldsymbol{\mu}}^{\text{exp}}, \hat{\boldsymbol{\Sigma}}^{\text{exp}})$$

$$= \mathcal{N} \left( \begin{bmatrix} \hat{\boldsymbol{\mu}}_1^{\text{exp}} \\ \hat{\boldsymbol{\mu}}_2^{\text{exp}} \\ \vdots \\ \hat{\boldsymbol{\mu}}_T^{\text{exp}} \end{bmatrix}, \begin{bmatrix} \hat{\boldsymbol{\Sigma}}_1^{\text{exp}} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \hat{\boldsymbol{\Sigma}}_2^{\text{exp}} & & \vdots \\ \vdots & & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \hat{\boldsymbol{\Sigma}}_T^{\text{exp}} \end{bmatrix} \right)$$

Estimated distribution



$p_{\text{exp}}(\boldsymbol{\tau})$

Structure is in line with Assumption 1.

\*Englert, P., Paraschos, A., Deisenroth, M. P., & Peters, J. (2013). Probabilistic model-based imitation learning. *Adaptive Behavior*, 21(5), 388–403.

# Forward Models

How to calculate  $p_{\theta}(\tau)$  ?


$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t$$

Forward models		Forward models	
Linear	$x_{t+1} \sim \mathcal{N}(A_t \tilde{x}_t + b_t, \Sigma_t)$	NN	$x_{t+1} \sim \mathcal{N}(\mu_{\theta_{\mu}}(\tilde{x}_t), \Sigma_{\theta_{\Sigma}}(\tilde{x}_t))$
Mixture of experts	$x_{t+1} \sim \sum_{k=1}^K \pi_k \mathcal{N}(A_{kt} \tilde{x}_t + b_{kt}, \Sigma_{kt})$	Koopman	$g(x_{t+1}) \sim \mathcal{N}(A_t g(x_t) + B_t u_t, \Sigma_t)$ $x_{t+1} = g^{-1}(g(x_{t+1}))$
RBF	$x_{t+1} = \sum_{k=1}^K \pi_k \mathcal{N}(\tilde{x}_t   \mu_k, \Sigma)$	More	?
GP	$x_{t+1} = f(\tilde{x}_t) + \epsilon$ $f \sim \text{GP}$ $\epsilon \sim \mathcal{N}(0, \Sigma_{\epsilon})$	$p(\mathbf{x}_{t+1} \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} \boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_{t+1})$ with $\boldsymbol{\mu}_{t+1} = \mathbb{E}_f[f(\mathbf{x}_t, \mathbf{u}_t)] = m_f(\mathbf{x}_t, \mathbf{u}_t)$ , $\boldsymbol{\Sigma}_{t+1} = \text{var}_f[f(\mathbf{x}_t, \mathbf{u}_t)] = \sigma_f^2(\mathbf{x}_t, \mathbf{u}_t)$ , $m_f(\tilde{\mathbf{x}}_{\star}) = \mathbf{k}_{\star}^{\top} \mathbf{K}^{-1} \mathbf{y}$ , $\sigma_f^2(\tilde{\mathbf{x}}_{\star}) = \mathbf{k}_{\star\star} - \mathbf{k}_{\star}^{\top} \mathbf{K}^{-1} \mathbf{k}_{\star}$	



# Control Policies

How to calculate  $p_{\theta}(\tau)$  ?



$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t$$

Control Policies		Control Policies	
Linear	$u_t \sim \mathcal{N}(A_t x_t + b_t, \Sigma_t)$	NN	$u_t \sim \mathcal{N}(\mu_{\theta_{\mu}}(x_t), \Sigma_{\theta_{\Sigma}}(x_t))$
Mixture of experts	$u_t \sim \sum_{k=1}^K \pi_k \mathcal{N}(A_{kt} x_t + b_{kt}, \Sigma_{kt})$	Koopman	$u_t \sim \mathcal{N}(A_t g(x_t), \Sigma_t)$ $x_t = g^{-1}(g(x_t))$
RBF	$u_t = \sum_{k=1}^K \pi_k \mathcal{N}(x_t   \mu_k, \Sigma)$	More	?
GP	$u_t = f(x_t) + \epsilon$ $f \sim \text{GP}$ $\epsilon \sim \mathcal{N}(0, \Sigma_{\epsilon})$		



# Calculating the trajectory distributions

How to calculate  $p_{\theta}(\tau)$  ?

Forward model      Control policy

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t$$

For which forward model and control policy combinations, do we have an analytical solution for this integral?

$$p(\mathbf{x}_{t+1}) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t \rightarrow \text{Same question?}$$

In PILCO and this paper\*, they found out a way to use GP forward models with nonlinear control policies with a lot of approximation which seem to work well.

$$\begin{aligned}x_{t+1} &\sim \mathcal{N}(A_t\tilde{x}_t + b_t, \Sigma_t) \\u_t &\sim \mathcal{N}(A_t x_t + b_t, \Sigma_t)\end{aligned}$$

Linear

$$\begin{aligned}g(x_{t+1}) &\sim \mathcal{N}(A_t g(x_t) + B_t u_t, \Sigma_t) \\u_t &\sim \mathcal{N}(A_t g(x_t), \Sigma_t) \\x_t &= g^{-1}(g(x_t))\end{aligned}$$

Koopman

[https://en.wikipedia.org/wiki/Conjugate\\_prior](https://en.wikipedia.org/wiki/Conjugate_prior)

\*Englert, P., Paraschos, A., Deisenroth, M. P., & Peters, J. (2013). Probabilistic model-based imitation learning. *Adaptive Behavior*, 21(5), 388–403.

# Calculating the trajectory distributions

How to calculate  $p_{\theta}(\tau)$  ?

In PILCO and this paper\*, they found out a way to use GP forward models with nonlinear control policies with a lot of approximation which seem to work well.

Forward model      Control policy

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t \quad \text{and} \quad p(\mathbf{x}_{t+1}) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t \quad \text{gives}$$

→ 
$$p(\mathbf{x}_{t+1}) = \int \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{u}_td\mathbf{x}_t$$

→ 
$$p(\mathbf{x}_{t+1}) = \int \int p(\mathbf{x}_{t+1}|\tilde{\mathbf{x}}_t)p(\tilde{\mathbf{x}}_t)df d\tilde{\mathbf{x}}_t \quad \text{with} \quad \begin{aligned} \mathbf{x}_{t+1} &= f(\tilde{\mathbf{x}}_t) + \epsilon \\ f &\sim \text{GP} \\ \epsilon &\sim \mathcal{N}(0, \Sigma_{\epsilon}) \end{aligned}$$

$$\begin{aligned} p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) &= \mathcal{N}(\mathbf{x}_{t+1}|\mu_{t+1}, \Sigma_{t+1}) & m_f(\tilde{\mathbf{x}}_{\star}) &= \mathbf{k}_{\star}^{\top} \mathbf{K}^{-1} \mathbf{y}, \\ \text{with } \mu_{t+1} &= \mathbb{E}_f[f(\mathbf{x}_t, \mathbf{u}_t)] = m_f(\mathbf{x}_t, \mathbf{u}_t), & \sigma_f^2(\tilde{\mathbf{x}}_{\star}) &= \mathbf{k}_{\star\star} - \mathbf{k}_{\star}^{\top} \mathbf{K}^{-1} \mathbf{k}_{\star} \\ \Sigma_{t+1} &= \text{var}_f[f(\mathbf{x}_t, \mathbf{u}_t)] = \sigma_f^2(\mathbf{x}_t, \mathbf{u}_t), \end{aligned}$$

\*Englert, P., Paraschos, A., Deisenroth, M. P., & Peters, J. (2013). Probabilistic model-based imitation learning. *Adaptive Behavior*, 21(5), 388–403.

# Calculating the trajectory distributions

How to calculate  $p_{\theta}(\tau)$  ?

In PILCO and this paper\*, they found out a way to use GP forward models with nonlinear control policies with a lot of approximation which seem to work well.

Forward model      Control policy

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t \quad \text{and} \quad p(\mathbf{x}_{t+1}) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t \quad \text{gives}$$

→ 
$$p(\mathbf{x}_{t+1}) = \int \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{u}_td\mathbf{x}_t$$

→ 
$$p(\mathbf{x}_{t+1}) = \int \int p(\mathbf{x}_{t+1}|\tilde{\mathbf{x}}_t)p(\tilde{\mathbf{x}}_t)df d\tilde{\mathbf{x}}_t$$

2. assumed Gaussian

3. Exact distribution not intractable. Approximate it by a Gaussian distribution.

Assumption 3

$$p(\mathbf{x}_{t+1}) \sim \mathcal{N}(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_{t+1})$$

Assumption 2

$$p(\tilde{\mathbf{x}}_t) \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\Sigma}}_t)$$

\*Englert, P., Paraschos, A., Deisenroth, M. P., & Peters, J. (2013). Probabilistic model-based imitation learning. *Adaptive Behavior*, 21(5), 388–403.

# Calculating the trajectory distributions for linear case

How to calculate  $p_{\theta}(\tau)$  ?

Forward model      Control policy

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)d\mathbf{u}_t \quad \text{and} \quad p(\mathbf{x}_{t+1}) = \int p(\mathbf{x}_{t+1}|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t \quad \text{gives}$$
$$p(\mathbf{x}_{t+1}) = \int \int p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)p(\mathbf{u}_t|\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{u}_td\mathbf{x}_t$$

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{A}_t\mathbf{x}_t + \mathbf{B}_t\mathbf{u}_t, \Sigma_t^d)$$

$$p(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(\mathbf{C}_t\mathbf{x}_t + \mathbf{d}_t, \Sigma_t^u)$$

HINT

$$\mathbb{E}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}[\mathbf{x}] + \mathbf{b}$$

$$\text{Var}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\text{Var}[\mathbf{x}]\mathbf{A}^T$$