



University of
BRISTOL

PhD Thesis: Bayesian Learning for Control in Multimodal Dynamical Systems

Aidan Scannell | Carl Henrik Ek | Arthur Richards

8th September 2022



Target



Mode 2



Start

Mode 1



Target



Mode 2



Start

Mode 1



Target



Mode 2



Start

Mode 1



Target



Mode 2



Start

Mode 1

Goals

Goal 1 Navigate to the target state x_f

Goal 2 Remain in the operable, desired dynamics mode k^*

Mode remaining navigation problem

$$\min_{\pi \in \Pi} \sum_{t=0}^T c(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) \quad (1a)$$

s.t. (1b)

(1c)

(1d)

(1e)

Mode remaining navigation problem

$$\min_{\pi \in \Pi} \sum_{t=0}^T c(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) \quad (1a)$$

s.t. (1b)

(1c)

$\mathbf{x}_0 = \mathbf{x}_0$ (1d)

(1e)

Mode remaining navigation problem

$$\min_{\pi \in \Pi} \sum_{t=0}^T c(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) \quad (1a)$$

s.t. (1b)

(1c)

$$\mathbf{x}_0 = \mathbf{x}_0 \quad (1d)$$

$$\mathbf{x}_T = \mathbf{x}_f \quad (1e)$$

Mode remaining navigation problem

$$\min_{\pi \in \Pi} \sum_{t=0}^T c(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) \quad (1a)$$

$$\text{s.t. } \mathbf{x}_{t+1} = f_k(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) + \boldsymbol{\epsilon}_k, \quad \text{if } \alpha(\mathbf{x}_t) = k \quad \forall t \in \{0, \dots, T-1\} \quad (1b)$$

(1c)

$$\mathbf{x}_0 = \mathbf{x}_0 \quad (1d)$$

$$\mathbf{x}_T = \mathbf{x}_f \quad (1e)$$

Mode remaining navigation problem

$$\min_{\pi \in \Pi} \sum_{t=0}^T c(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) \quad (1a)$$

$$\text{s.t. } \mathbf{x}_{t+1} = f_k(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) + \boldsymbol{\epsilon}_k, \quad \text{if } \alpha(\mathbf{x}_t) = k \quad \forall t \in \{0, \dots, T-1\} \quad (1b)$$

$$\alpha(\mathbf{x}_t) = k^* \quad \forall t \in \{0, \dots, T-1\} \quad (1c)$$

$$\mathbf{x}_0 = \mathbf{x}_0 \quad (1d)$$

$$\mathbf{x}_T = \mathbf{x}_f \quad (1e)$$

But dynamics are not known a priori . . .

$$\min_{\pi \in \Pi} \sum_{t=0}^T c(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) \quad (2a)$$

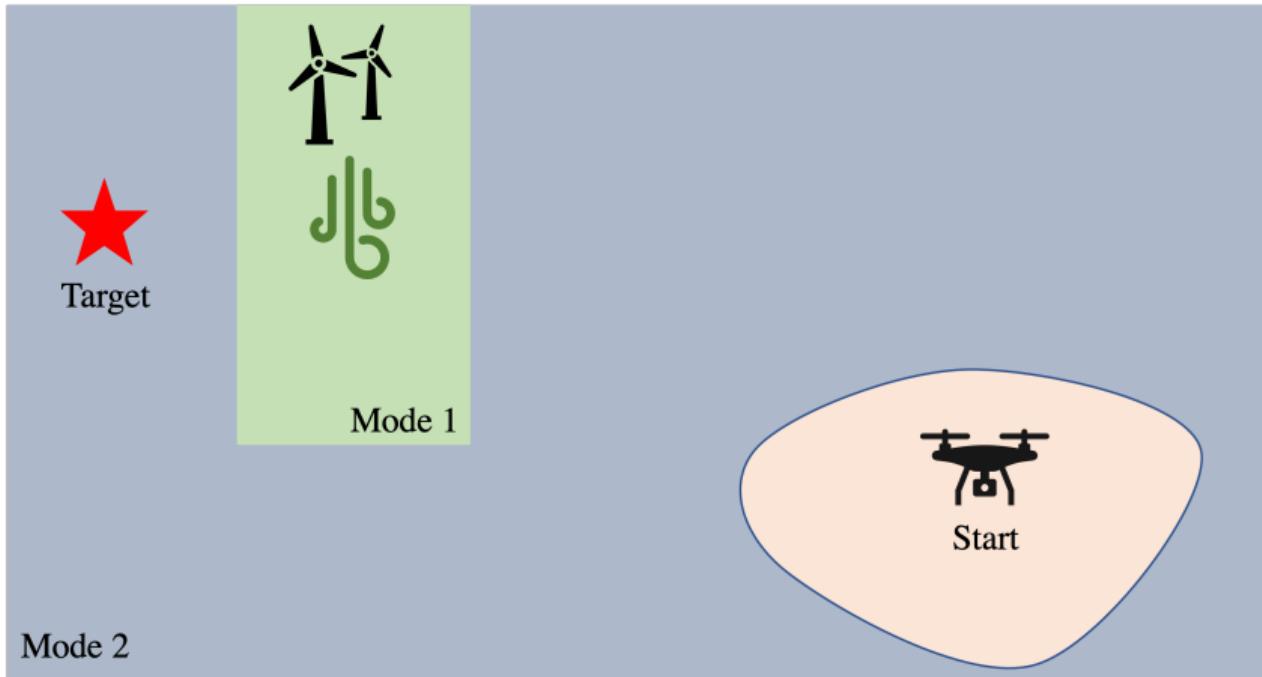
$$\text{s.t. } \mathbf{x}_{t+1} = \mathbf{f}_k(\mathbf{x}_t, \pi(\mathbf{x}_t, t)) + \boldsymbol{\epsilon}_k, \quad \text{if } \alpha(\mathbf{x}_t) = k \quad \forall t \in \{0, \dots, T-1\} \quad (2b)$$

$$\alpha(\mathbf{x}_t) = k^* \quad \forall t \in \{0, \dots, T-1\} \quad (2c)$$

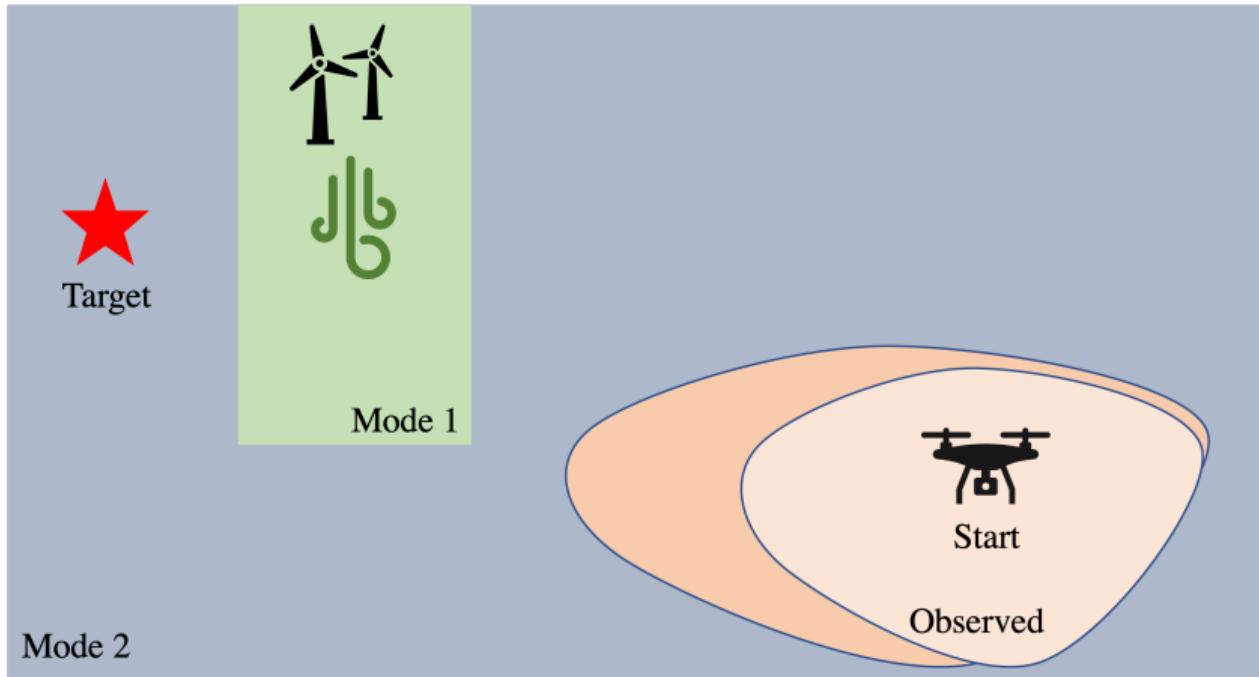
$$\mathbf{x}_0 = \mathbf{x}_0 \quad (2d)$$

$$\mathbf{x}_T = \mathbf{x}_f \quad (2e)$$

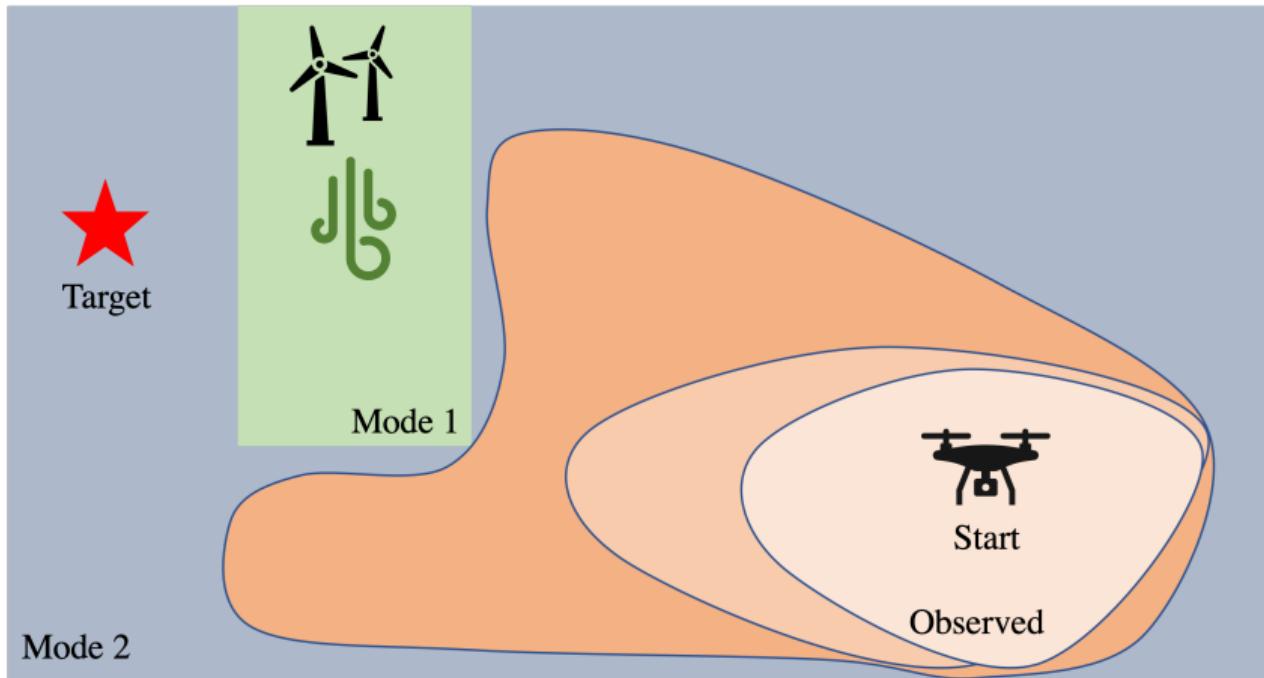
Solve using model-based reinforcement learning?



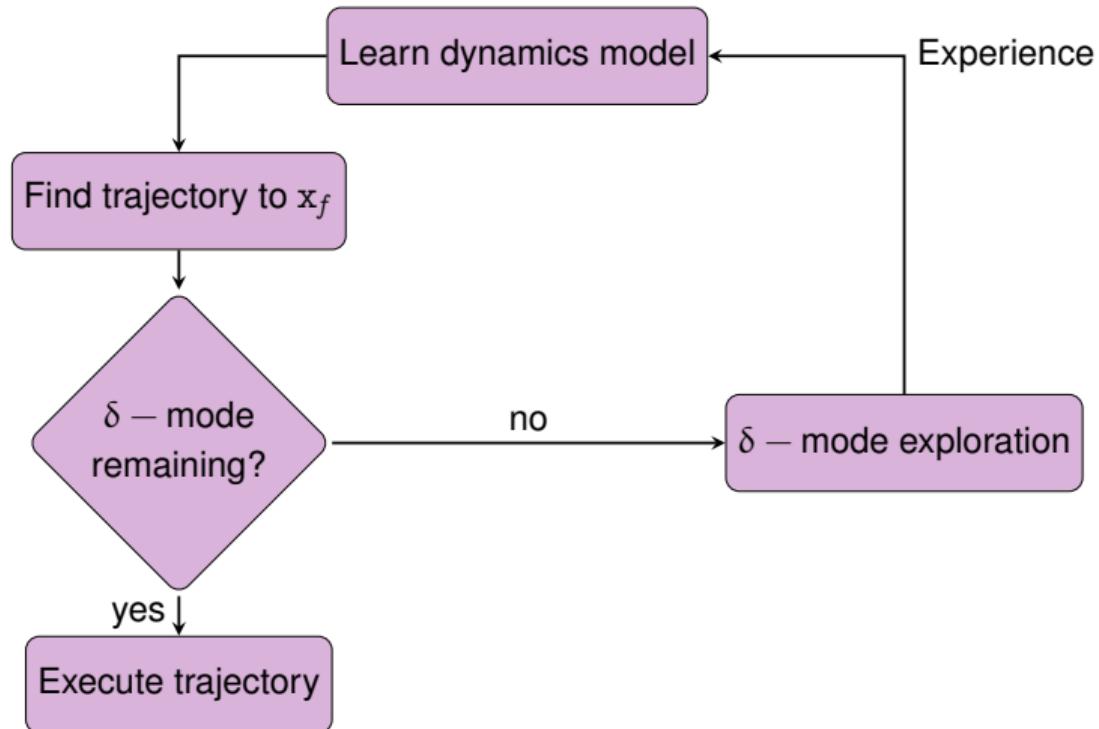
Solve using model-based reinforcement learning?



Solve using model-based reinforcement learning?



Solve using model-based reinforcement learning?



Contributions

1. Model learning

Contributions

1. Model learning
2. Mode remaining trajectory optimisation

Contributions

1. Model learning
2. Mode remaining trajectory optimisation
 - ▶ via latent geometry

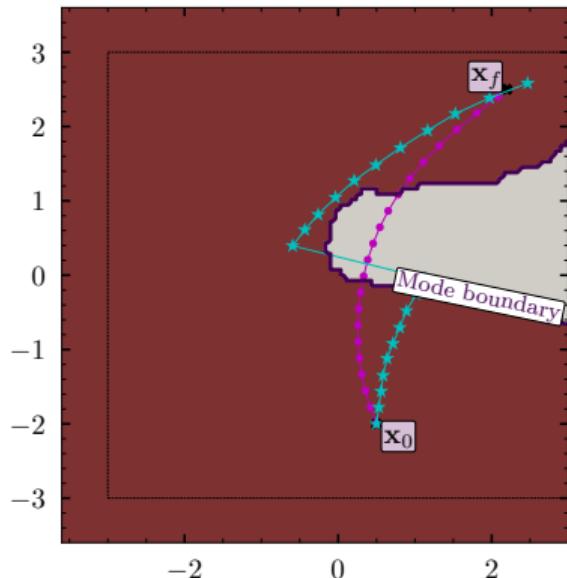
Contributions

1. Model learning
2. Mode remaining trajectory optimisation
 - ▶ via latent geometry
 - ▶ control as inference

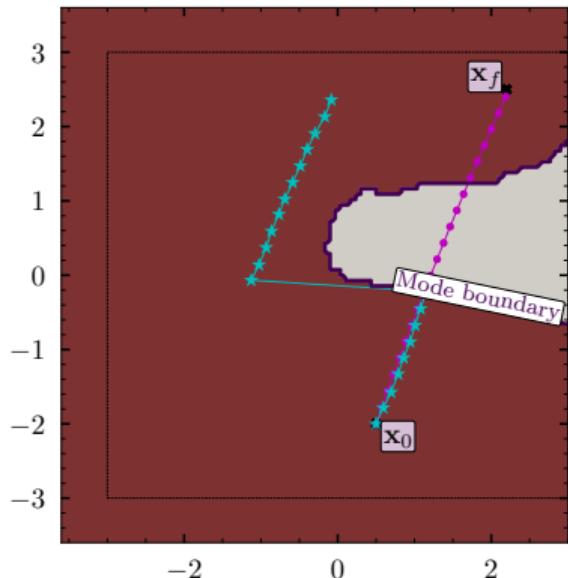
Contributions

1. Model learning
2. Mode remaining trajectory optimisation
 - ▶ via latent geometry
 - ▶ control as inference
3. Mode remaining exploration for model-based reinforcement learning

Model learning - Gaussian processes don't work...



—●— Dynamics —★— Environment



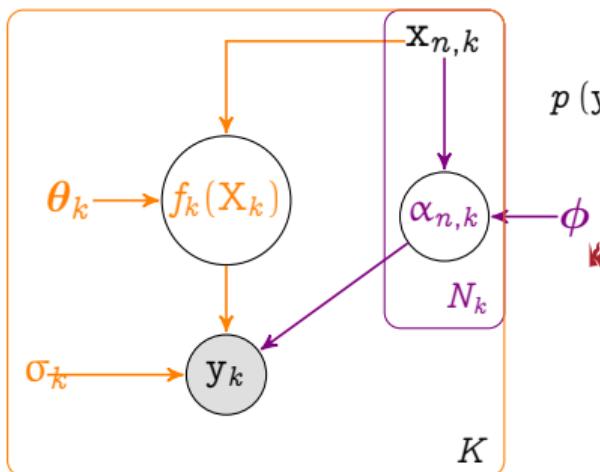
—●— Dynamics —★— Environment

Model learning - mixture models?

MoE marginal likelihood

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \prod_{n=1}^N \sum_{k=1}^K \underbrace{\Pr(\alpha_n = k | \mathbf{x}_n, \boldsymbol{\phi})}_{\text{gating network}} \underbrace{p(y_n | \alpha_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)}_{\text{expert } k}, \quad (3)$$

Model learning - mixtures of nonparametric experts



$$p(y | X, \theta, \phi) = \sum_{\alpha} p(\alpha | X, \phi) \left[\prod_{k=1}^K p(\{y_n : \alpha_n = k\} | \{x_n : \alpha_n = k\}, \theta_k) \right]$$

☞ Sum over exponentially many (K^N) sets of assignments

► $\alpha = \{\alpha_1, \dots, \alpha_N\}$

Model learning - Parameterise the nonparametric model?

- » Like a sparse GP parameterises a GP...

Model learning - Parameterise the nonparametric model?

- » Like a sparse GP parameterises a GP...
- » GP prior where $X_k = \{x_n : \alpha_t = k\}$

$$f_k(X_k) \sim \mathcal{N}(\mu_k(X_k), k_k(X_k, X_k))$$

Model learning - Parameterise the nonparametric model?

- » Like a sparse GP parameterises a GP...
- » GP prior where $X_k = \{x_n : \alpha_t = k\}$

$$f_k(X_k) \sim \mathcal{N}(\mu_k(X_k), k_k(X_k, X_k))$$

- » Augment with inducing points

$$f_k(\zeta_k) \sim \mathcal{N}(\mu_k(\zeta_k), k_k(\zeta_k, \zeta_k))$$

Model learning - Parameterise the nonparametric model?

- Like a sparse GP parameterises a GP...
- GP prior where $X_k = \{x_n : \alpha_t = k\}$

$$f_k(X_k) \sim \mathcal{N}(\mu_k(X_k), k_k(X_k, X_k))$$

- Augment with inducing points

$$f_k(\zeta_k) \sim \mathcal{N}(\mu_k(\zeta_k), k_k(\zeta_k, \zeta_k))$$

- FITC for MoGPE?

$$p(y | f(\zeta)) \approx \prod_{n=1}^N p(y_n | f(\zeta)) = \prod_{n=1}^N \sum_{k=1}^K \Pr(\alpha_n = k | x_n, \phi) \prod_{k=1}^K p(y_n | f_k(\zeta_k)).$$

Model learning - Parameterise the nonparametric model?

- Like a sparse GP parameterises a GP...
- GP prior where $X_k = \{x_n : \alpha_t = k\}$

$$f_k(X_k) \sim \mathcal{N}(\mu_k(X_k), k_k(X_k, X_k))$$

- Augment with inducing points

$$f_k(\zeta_k) \sim \mathcal{N}(\mu_k(\zeta_k), k_k(\zeta_k, \zeta_k))$$

- FITC for MoGPE?

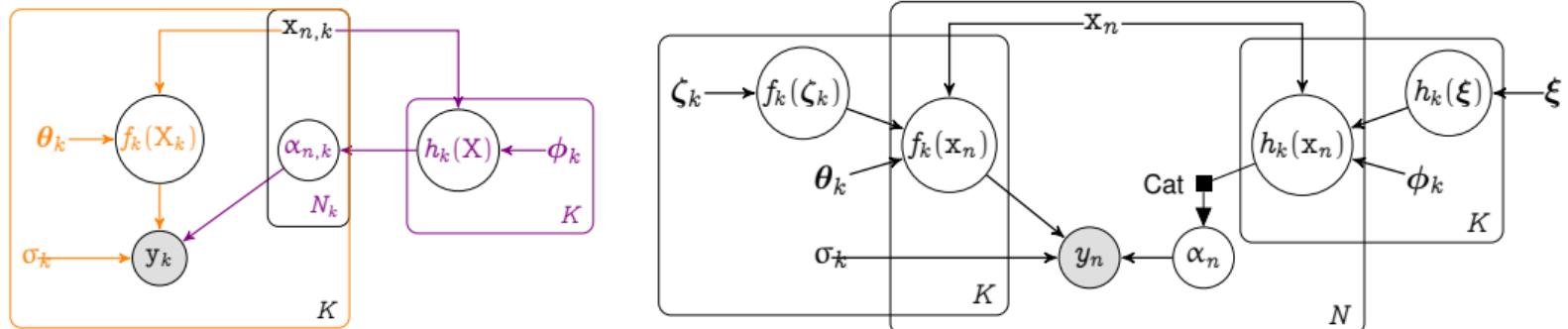
$$p(y | f(\zeta)) \approx \prod_{n=1}^N p(y_n | f(\zeta)) = \prod_{n=1}^N \sum_{k=1}^K \Pr(\alpha_n = k | x_n, \phi) \prod_{k=1}^K p(y_n | f_k(\zeta_k)).$$

- Assumes inducing variables $\{f_k(\zeta_k)\}_{k=1}^K$, are a sufficient statistic for latent function values $\{f_k(X_k)\}_{k=1}^K$ AND the set of assignments α .

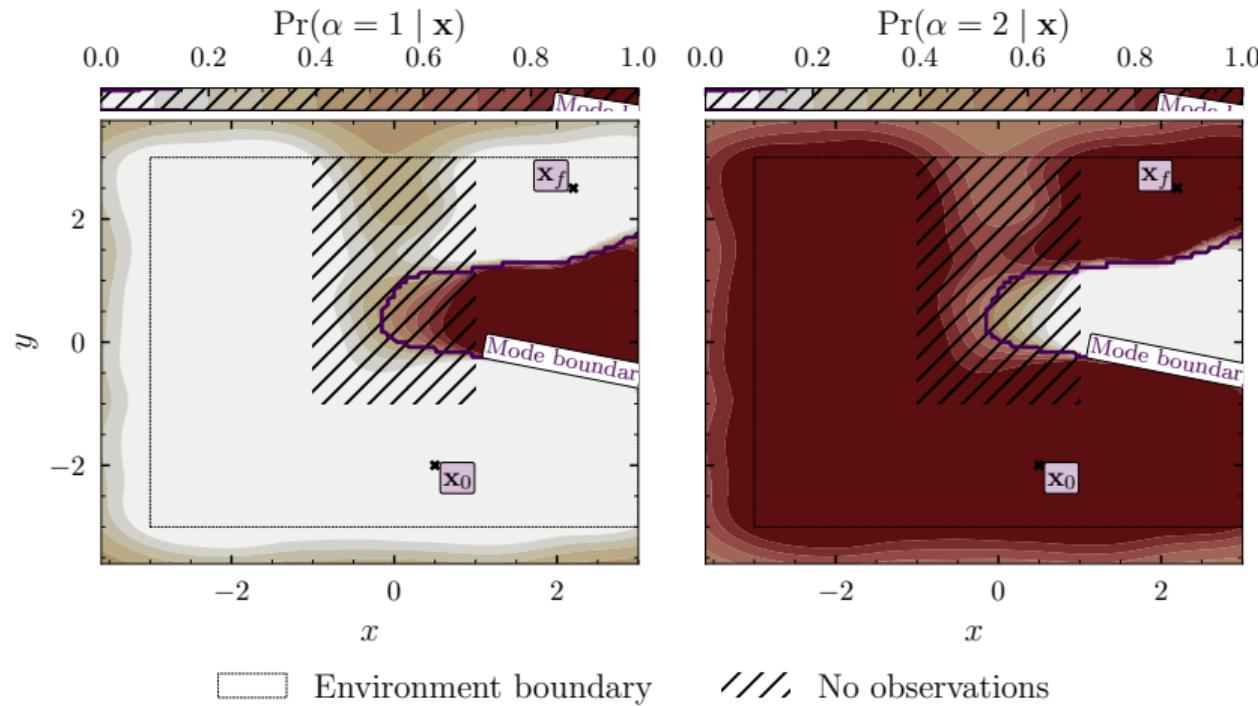
Model learning - Parameterise the nonparametric model?

☛ Approximate marginal likelihood

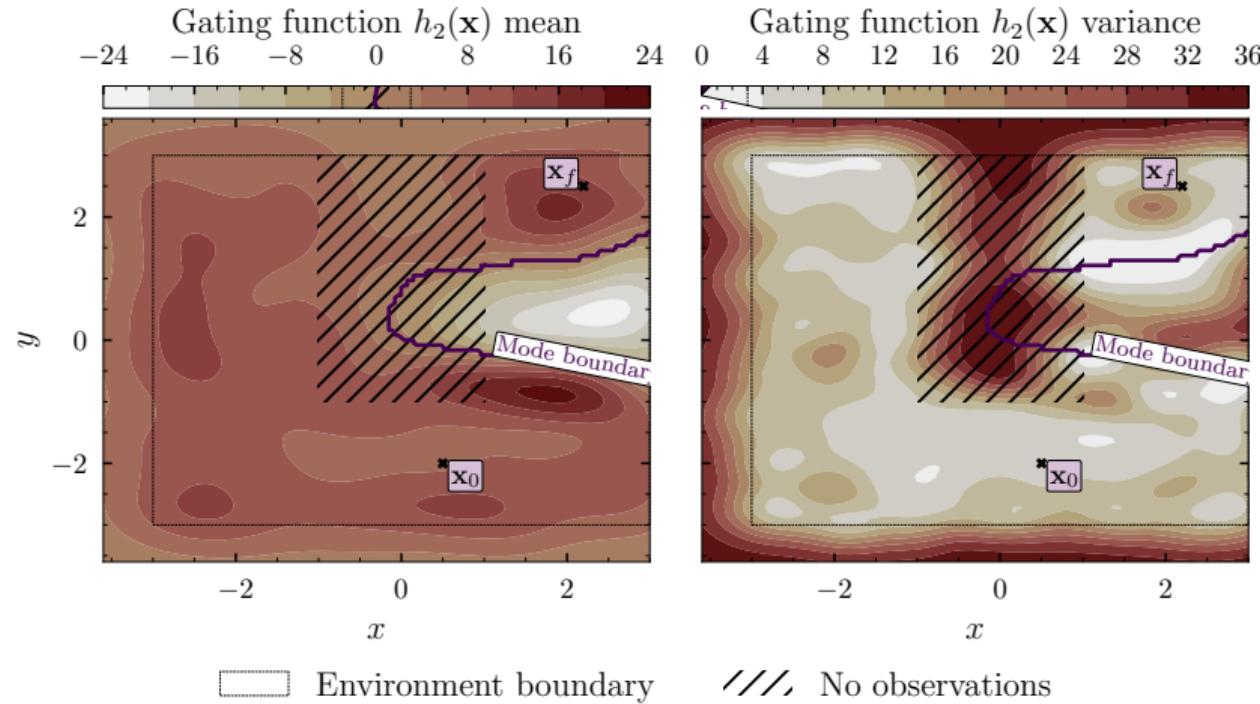
$$p(\mathbf{y} | \mathbf{X}) \approx \mathbb{E}_{p(h(\xi))p(f(\zeta))} \left[\prod_{n=1}^N \sum_{k=1}^K \Pr(\alpha_n = k | h(\xi)) p(y_n | f_k(\zeta_k)) \right] \quad (4)$$



Model learning - latent spaces for planning



Model learning - latent spaces for planning



Mode remaining control

☛ Goals

Mode remaining control

❖ Goals

- ▶ Navigate to the target state x_f

Mode remaining control

Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

Mode remaining control

❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

❖ Assumptions

Mode remaining control

�� Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

�� Assumptions

- ▶ Desired dynamics mode is *known a priori*

Mode remaining control

❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

❖ Assumptions

- ▶ Desired dynamics mode is *known a priori*
- ▶ Prior access to environment

Mode remaining control

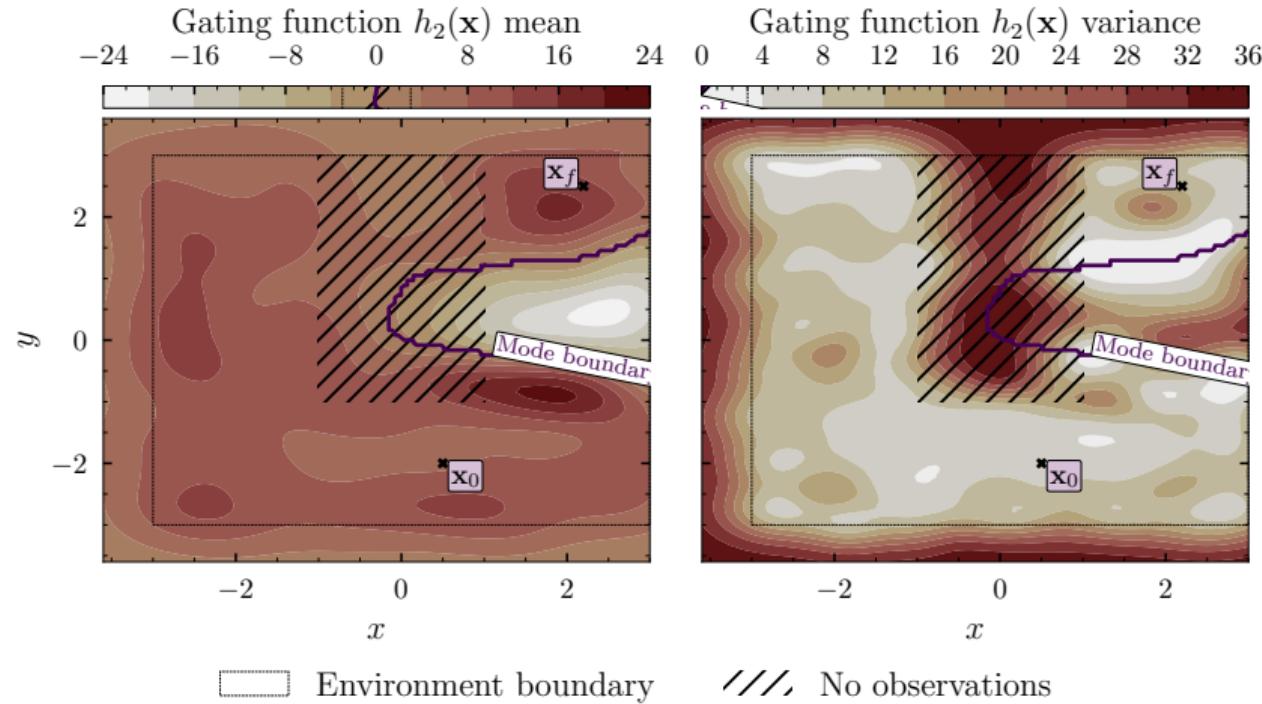
❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

❖ Assumptions

- ▶ Desired dynamics mode is *known a priori*
- ▶ Prior access to environment
 - ▶ Such that a state transition data set \mathcal{D} has been collected

Mode remaining control - via latent geometry



Mode remaining control - via latent geometry

- ❖ Desired mode's gating function $h_{k^*} : \mathcal{X} \rightarrow \mathbb{R}$

Mode remaining control - via latent geometry

- ❖ Desired mode's gating function $h_{k^*} : \mathcal{X} \rightarrow \mathbb{R}$
- ❖ State trajectory $\bar{x} : [t_0, t_f] \rightarrow \mathcal{X}$

Mode remaining control - via latent geometry

- Desired mode's gating function $h_{k^*} : \mathcal{X} \rightarrow \mathbb{R}$
- State trajectory $\bar{x} : [t_0, t_f] \rightarrow \mathcal{X}$
- Length minimising trajectories encode mode remaining behaviour

$$\min \text{Length}(h_{k^*}(\bar{x})) = \min \int_{t_0}^{t_f} \|\dot{x}(t)\|_{G(x(t))} dt \quad (5)$$

where,

$$\|\dot{x}(t)\|_{G(x(t))} = \sqrt{\dot{x}(t) G_{x_t} \dot{x}(t)} \quad (6)$$

Mode remaining control - via latent geometry

- Desired mode's gating function $h_{k^*} : \mathcal{X} \rightarrow \mathbb{R}$
- State trajectory $\bar{x} : [t_0, t_f] \rightarrow \mathcal{X}$
- Length minimising trajectories encode mode remaining behaviour

$$\min \text{Length}(h_{k^*}(\bar{x})) = \min \int_{t_0}^{t_f} \|\dot{x}(t)\|_{G(x(t))} dt \quad (5)$$

where,

$$\|\dot{x}(t)\|_{G(x(t))} = \sqrt{\dot{x}(t) G_{x_t} \dot{x}(t)} \quad (6)$$

- But ignores *epistemic uncertainty*...

Mode remaining control - via latent geometry

- Metric depends on Jacobian¹

$$G_{x_t} = J_{x_t}^T J_{x_t} \quad (7)$$

[3] Tosi et al. "Metrics for Probabilistic Geometries". 2014.

Mode remaining control - via latent geometry

- Metric depends on Jacobian¹

$$G_{x_t} = J_{x_t}^T J_{x_t} \quad (7)$$

- which is Normally distributed

$$J \sim \mathcal{N}(\mu_J, \Sigma_J) \quad (8)$$

[3] Tosi et al. "Metrics for Probabilistic Geometries". 2014.

Mode remaining control - via latent geometry

- Metric depends on Jacobian¹

$$G_{x_t} = J_{x_t}^T J_{x_t} \quad (7)$$

- which is Normally distributed

$$J \sim \mathcal{N}(\mu_J, \Sigma_J) \quad (8)$$

- so metric follows non-central Wishart distribution

$$G \sim W_D(p, \Sigma_J, \mathbb{E}[J^T] \mathbb{E}[J]) \quad (9)$$

[3] Tosi et al. "Metrics for Probabilistic Geometries". 2014.

Mode remaining control - via latent geometry

- Metric depends on Jacobian¹

$$G_{x_t} = J_{x_t}^T J_{x_t} \quad (7)$$

- which is Normally distributed

$$J \sim \mathcal{N}(\mu_J, \Sigma_J) \quad (8)$$

- so metric follows non-central Wishart distribution

$$G \sim \mathcal{W}_D(p, \Sigma_J, \mathbb{E}[J^T] \mathbb{E}[J]) \quad (9)$$

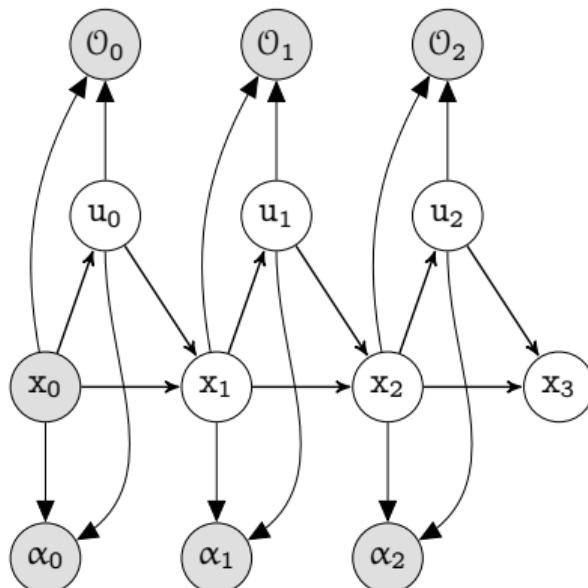
- Expected metric increases length of trajectories in regions of high *epistemic uncertainty*

$$\mathbb{E}[G] = \mathbb{E}[J^T] \mathbb{E}[J] + \lambda \Sigma_J \quad (10)$$

[3] Tosi et al. "Metrics for Probabilistic Geometries". 2014.

Mode remaining control as probabilistic inference

$$\Pr(\mathcal{O}_t = 1 \mid \mathbf{x}_t, \mathbf{u}_t) \propto \exp(-\gamma c(\mathbf{x}_t, \mathbf{u}_t))$$



Mode remaining control as probabilistic inference

Goal $p(\bar{x}, \bar{u} | x_0, \mathcal{O}_{0:T} = 1, \alpha_{0:T} = k^*)$

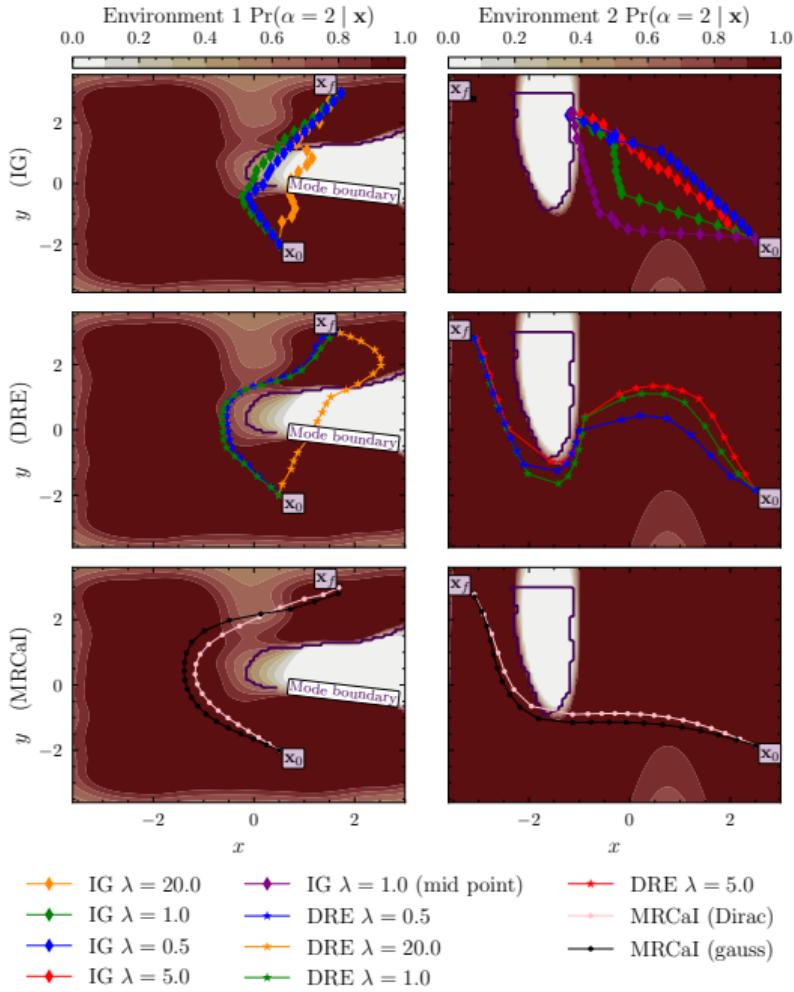
Mode remaining control as probabilistic inference

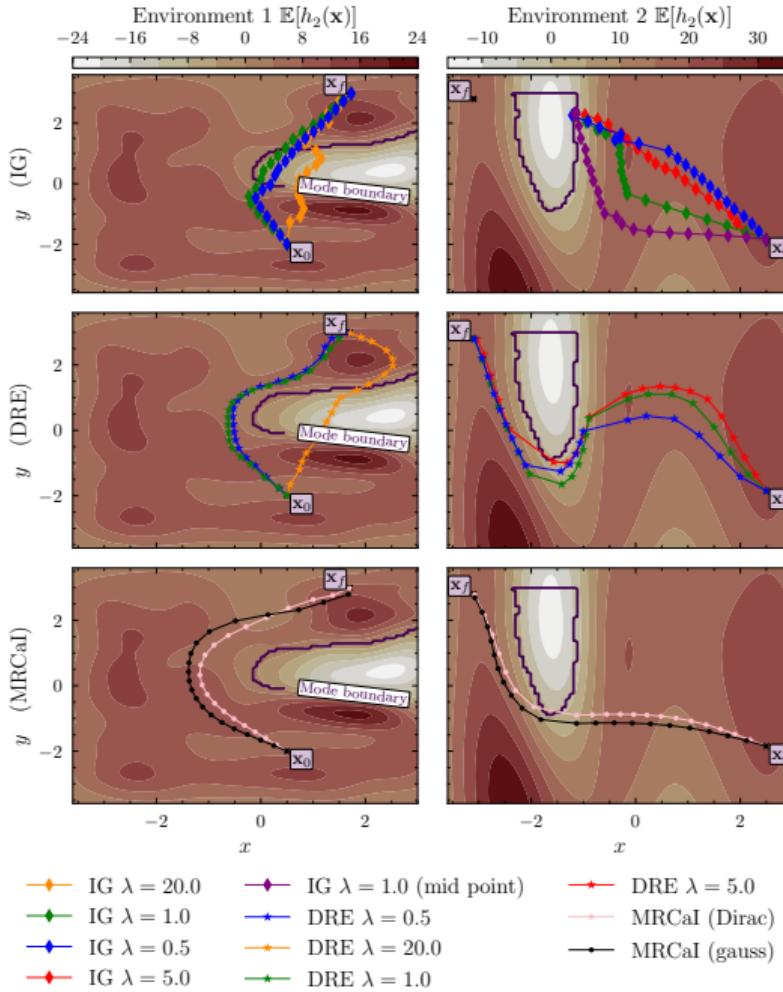
- Goal $p(\bar{x}, \bar{u} | x_0, \mathcal{O}_{0:T} = 1, \alpha_{0:T} = k^*)$
- Variational inference (lower bound $p(\mathcal{O}_{0:T} = 1, \alpha_{0:T} = k^* | x_0)$)

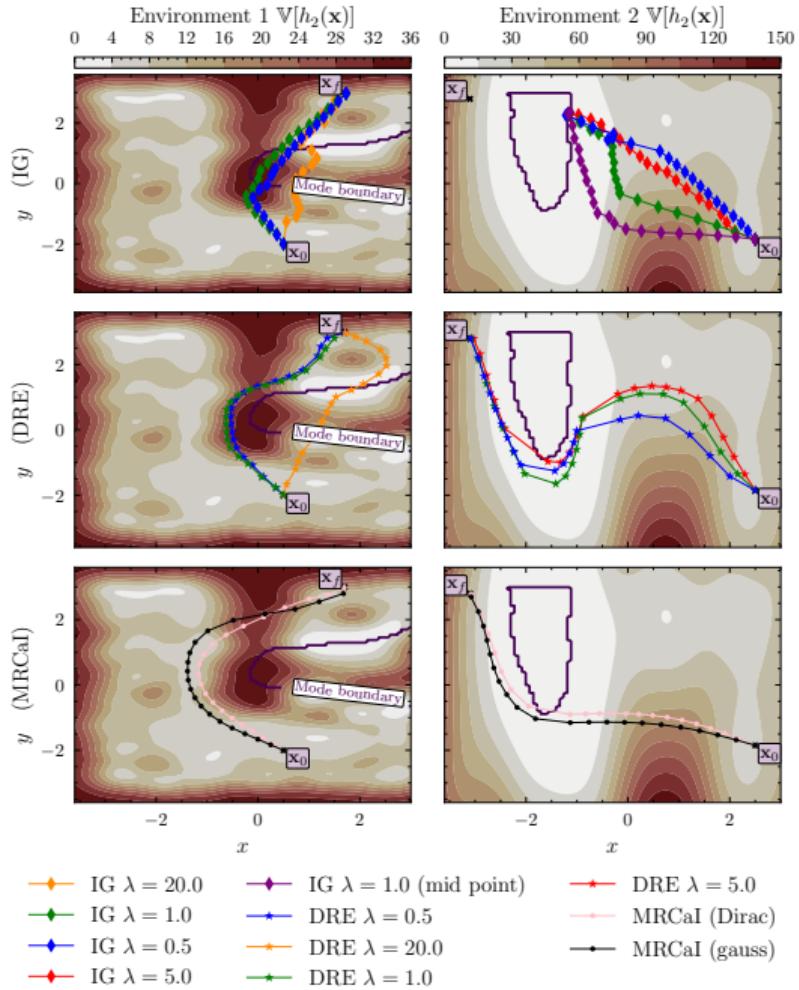
$$\mathcal{L}_{\text{mode}} = - \sum_{t=0}^T \underbrace{\mathbb{E}_{q(x_t|x_0, \alpha_{0:T}=k_{0:t-1}^*) q(u_t)} [c(x_t, u_t)]}_{\text{expected cost}} \quad (11)$$

$$+ \sum_{t=0}^T \underbrace{\mathbb{E}_{q(x_t|x_0, \alpha_{0:T}=k_{0:t-1}^*)} [\log \Pr(\alpha_t = k^* | x_t)]}_{\text{mode remaining term}}$$

$$+ \sum_{t=0}^{T-1} \underbrace{H[u_t]}_{\text{entropy}} \quad (12)$$







Mode remaining exploration for MBRL

❖ Goals

Mode remaining exploration for MBRL

❖ Goals

- ▶ Navigate to the target state x_f

Mode remaining exploration for MBRL

❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

Mode remaining exploration for MBRL

❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

❖ Assumptions

Mode remaining exploration for MBRL

❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

❖ Assumptions

- ▶ Desired dynamics mode is *known a priori*

Mode remaining exploration for MBRL

❖ Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

❖ Assumptions

- ▶ Desired dynamics mode is *known a priori*
- ▶ No access to environment a priori

Mode remaining exploration for MBRL

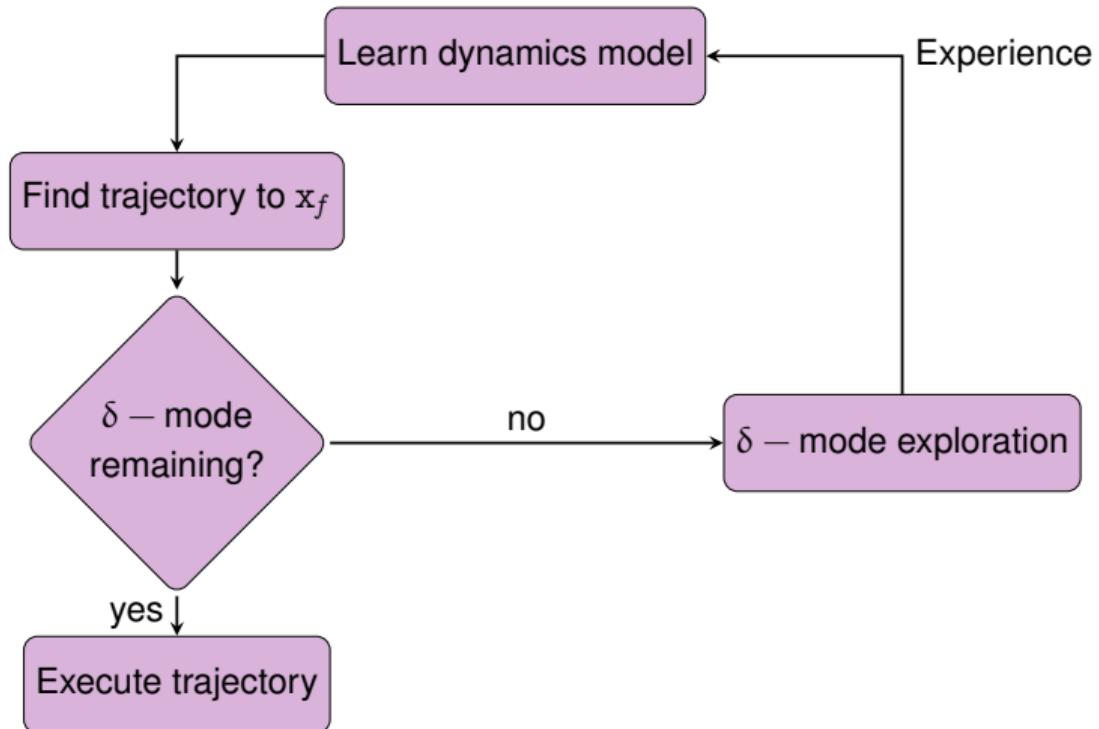
Goals

- ▶ Navigate to the target state x_f
- ▶ Remain in desired dynamics mode

Assumptions

- ▶ Desired dynamics mode is *known a priori*
- ▶ No access to environment a priori
 - ▶ Only a local state transition data set around start state

Mode remaining exploration for MBRL - the loop



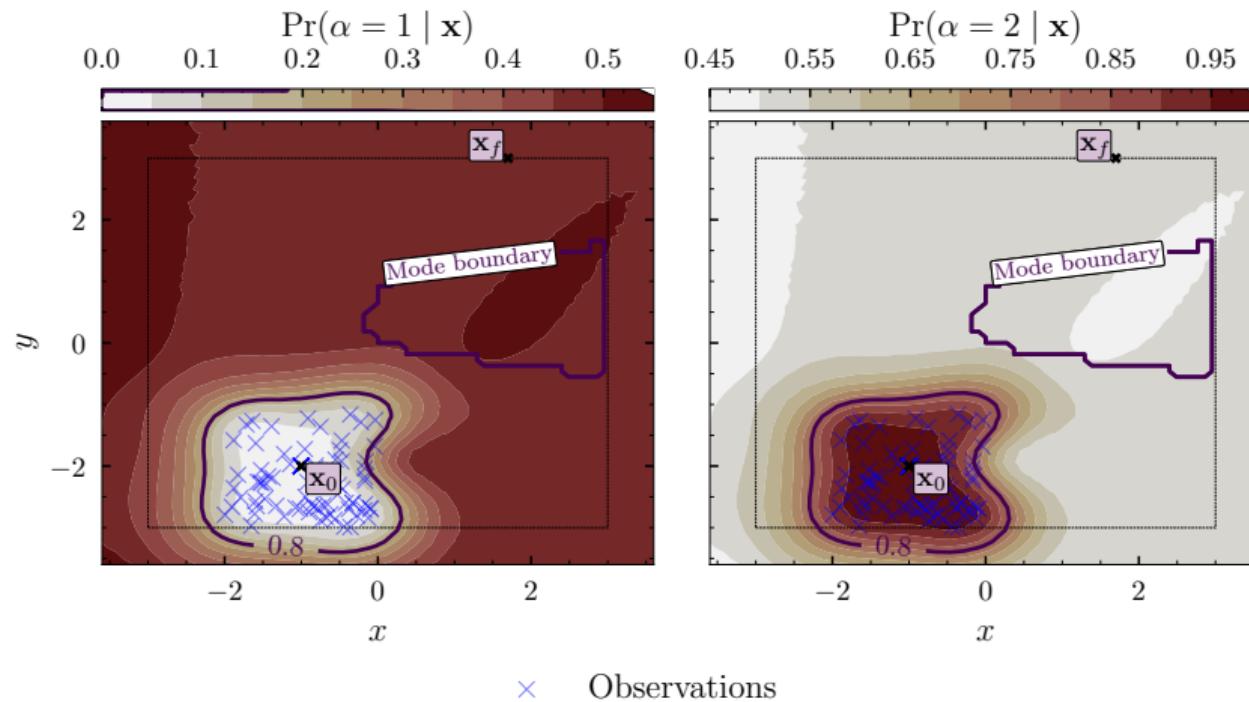
Mode remaining exploration for MBRL - information-based objective

$$\max_{\pi \in \Pi} \underbrace{\mathcal{H}[h_{k^*}(\bar{x}) | \bar{x}, \mathcal{D}_{0:i-1}]}_{\text{joint gating entropy}} + \sum_{t=1}^{T-1} \mathbb{E} \left[\underbrace{-(x_t - x_f)^T Q (x_t - x_f)}_{\text{state difference}} - \underbrace{u_t^T R u_t}_{\text{control cost}} \right] \quad (13a)$$

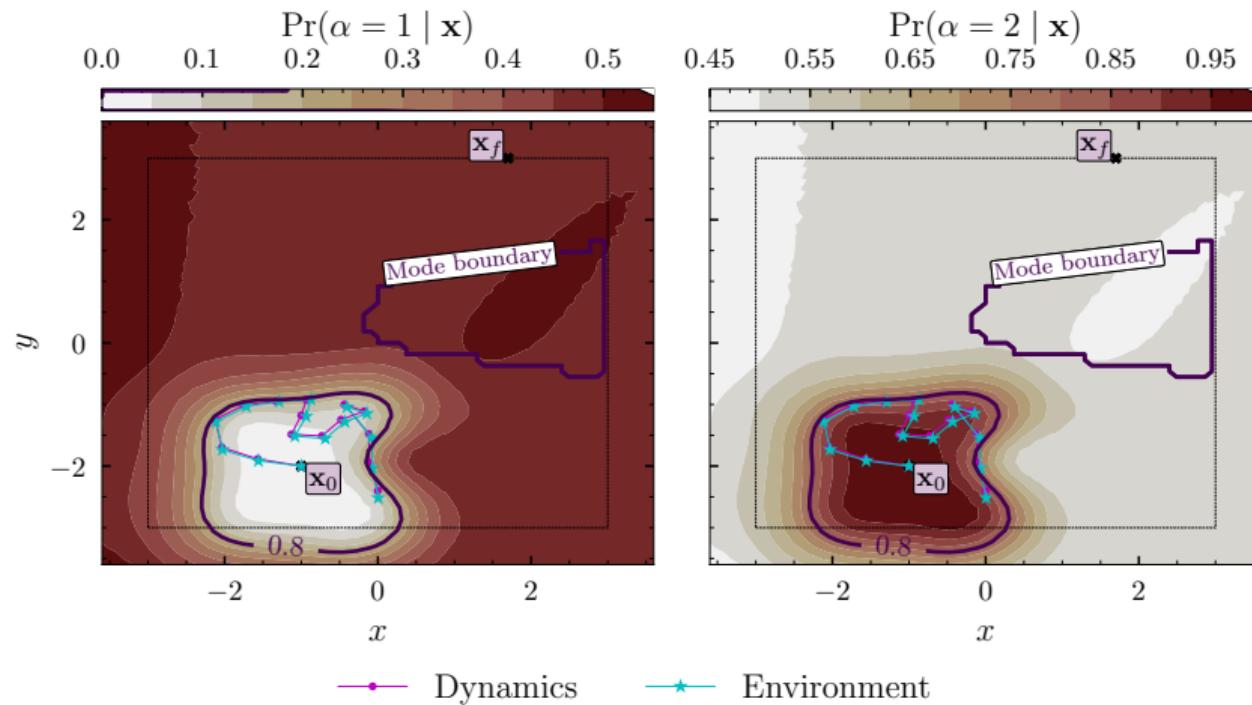
Mode remaining exploration for MBRL - chance constraints

$$\Pr(\alpha_t = k^* \mid \mathbf{x}_0, \mathbf{u}_{0:t}, \boldsymbol{\alpha}_{0:t-1} = k^*) \geq 1 - \delta \quad \forall t \in \{0, \dots, T\}$$

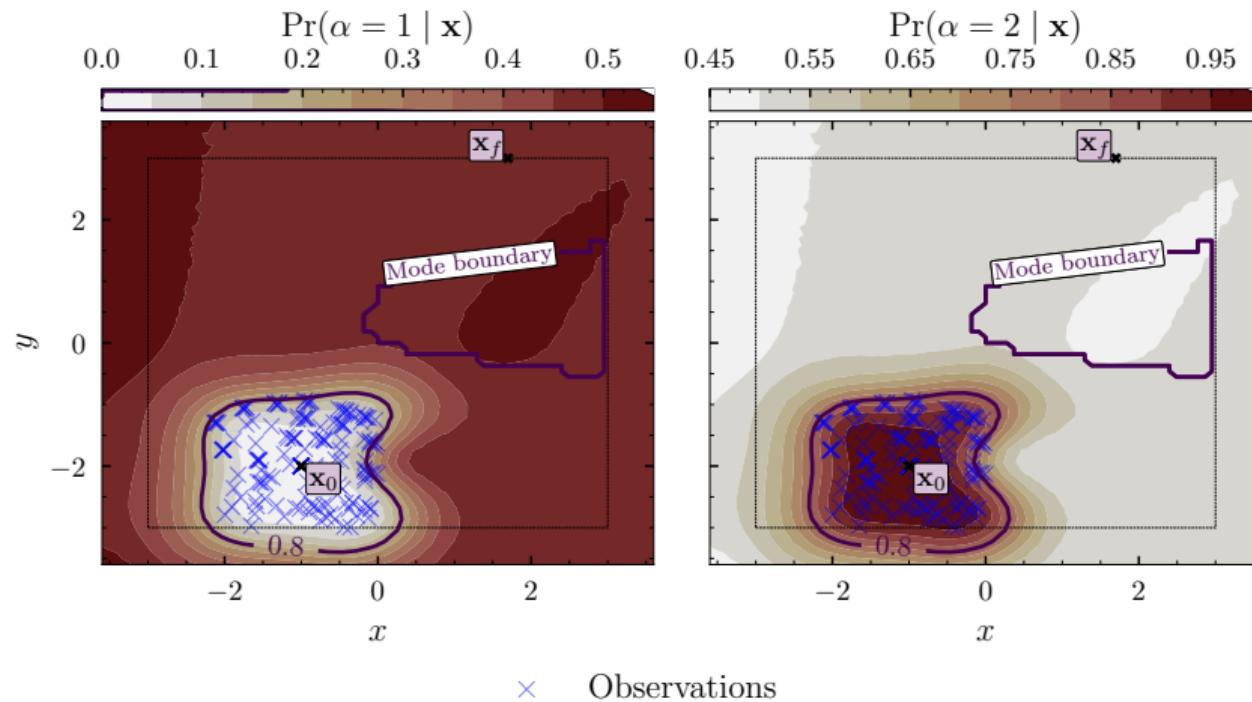
Mode remaining exploration for MBRL - iteration 0



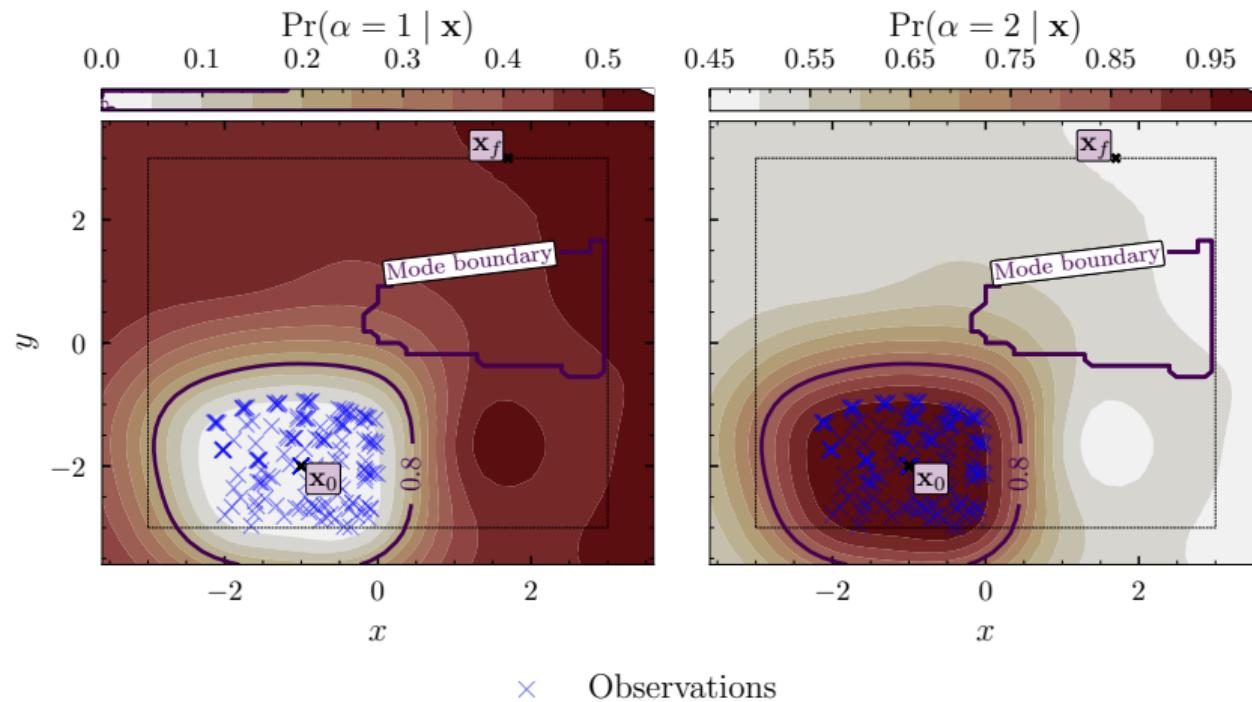
Mode remaining exploration for MBRL - iteration 0



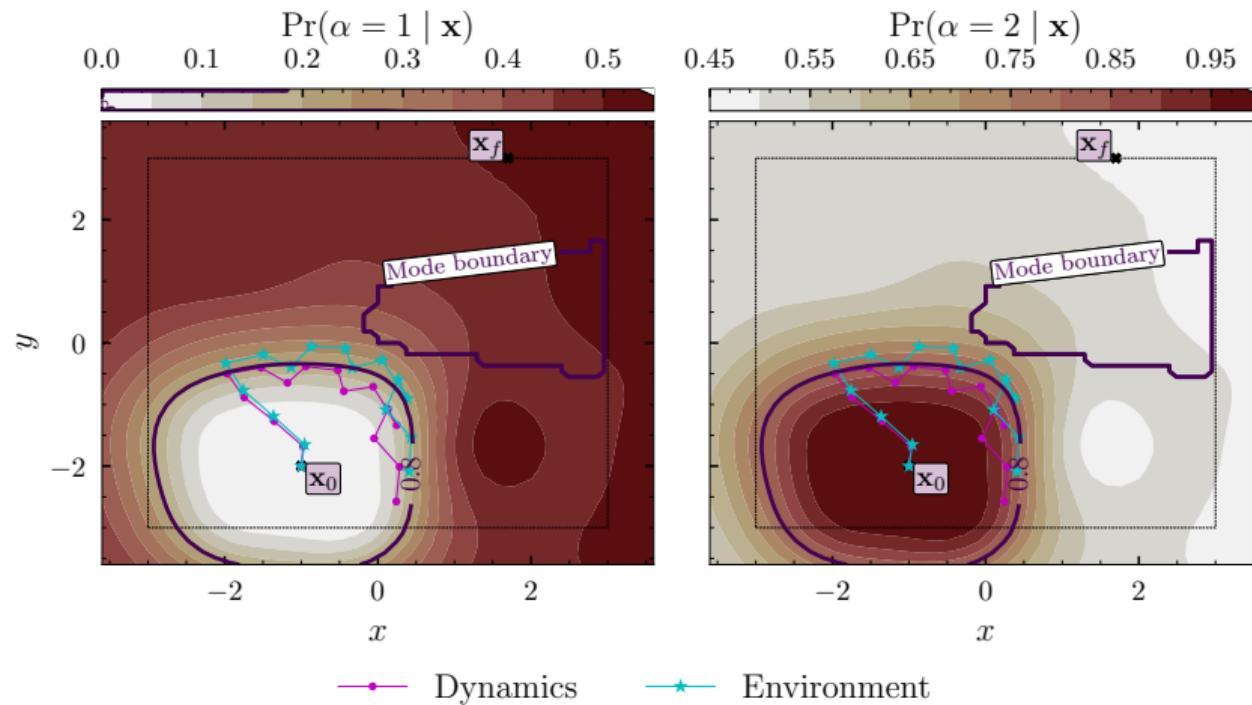
Mode remaining exploration for MBRL - iteration 0



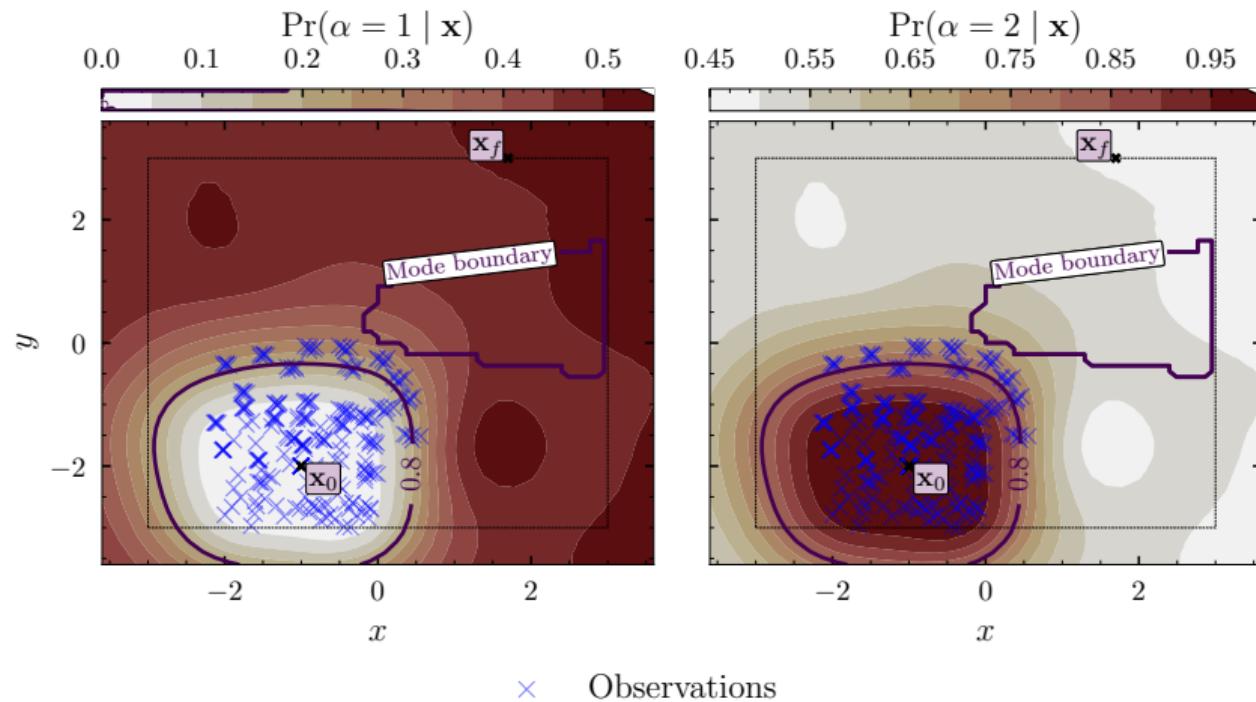
Mode remaining exploration for MBRL - iteration 1



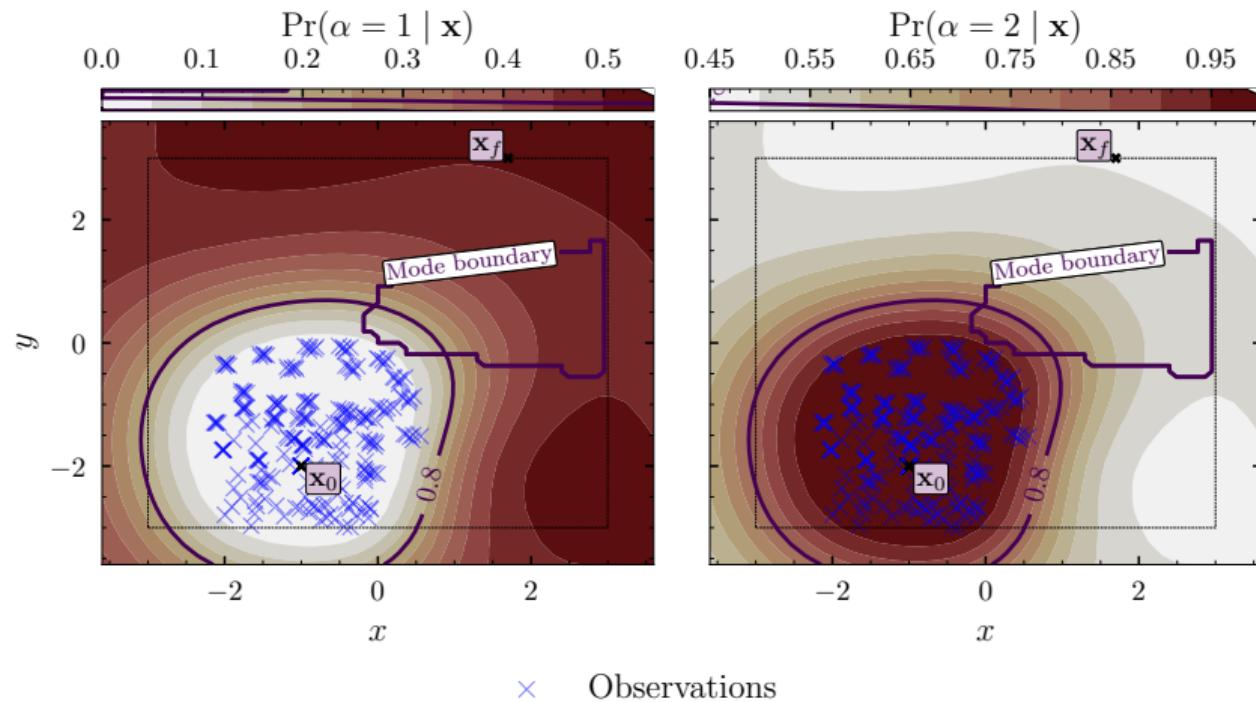
Mode remaining exploration for MBRL - iteration 1



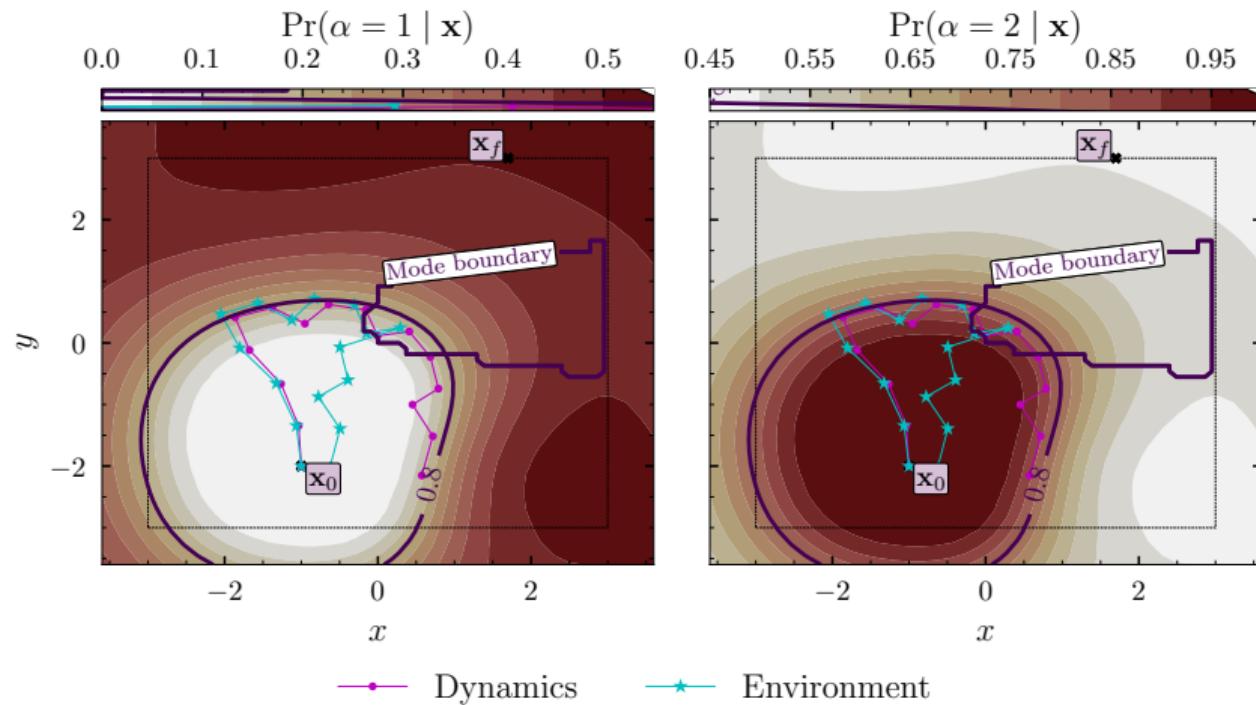
Mode remaining exploration for MBRL - iteration 1



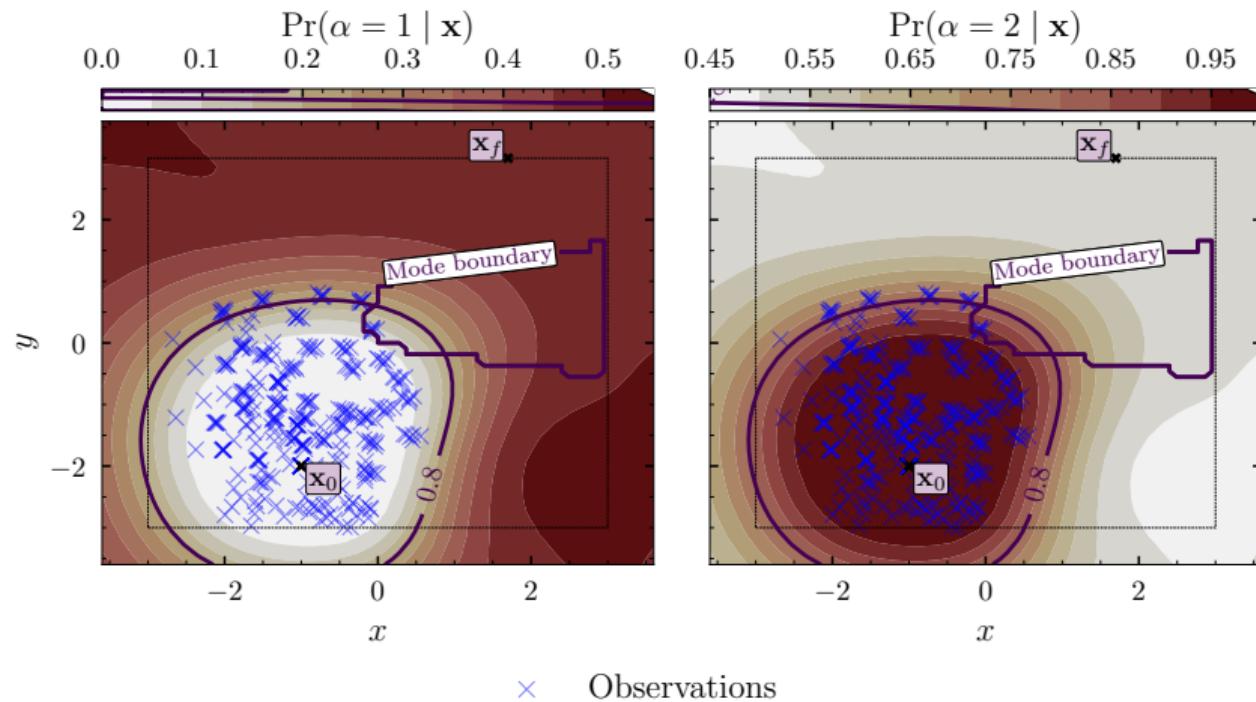
Mode remaining exploration for MBRL - iteration 2



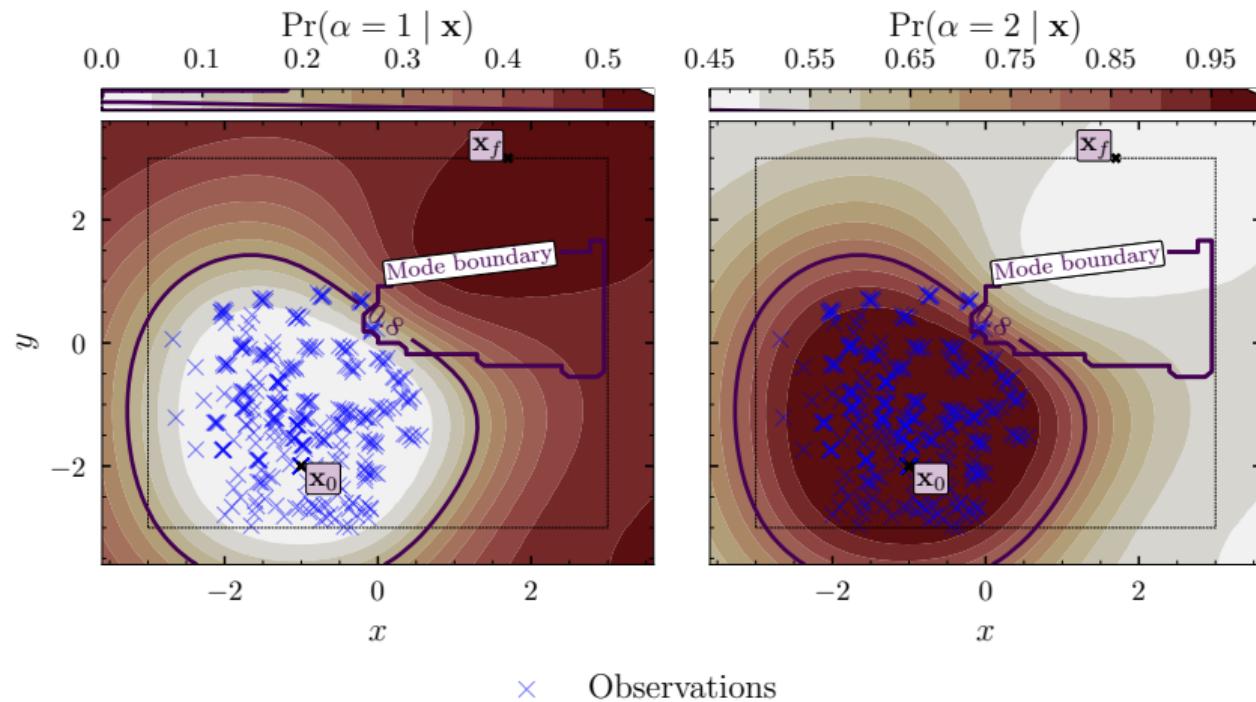
Mode remaining exploration for MBRL - iteration 2



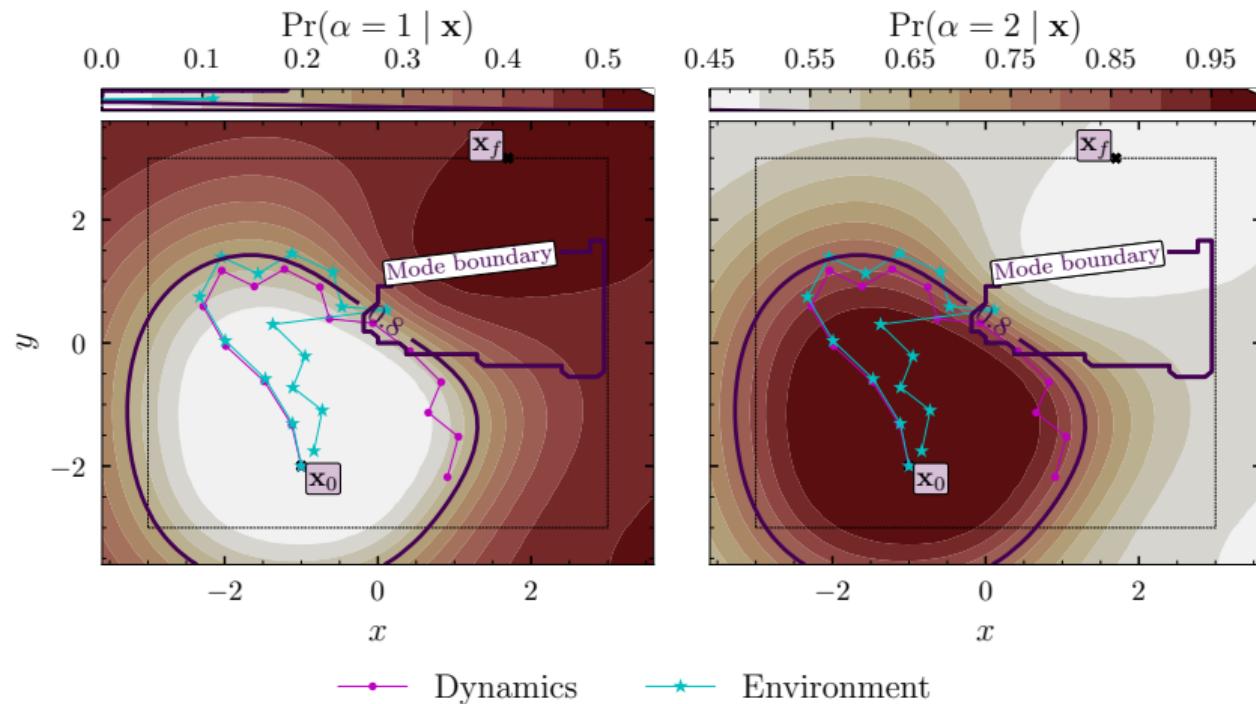
Mode remaining exploration for MBRL - iteration 2



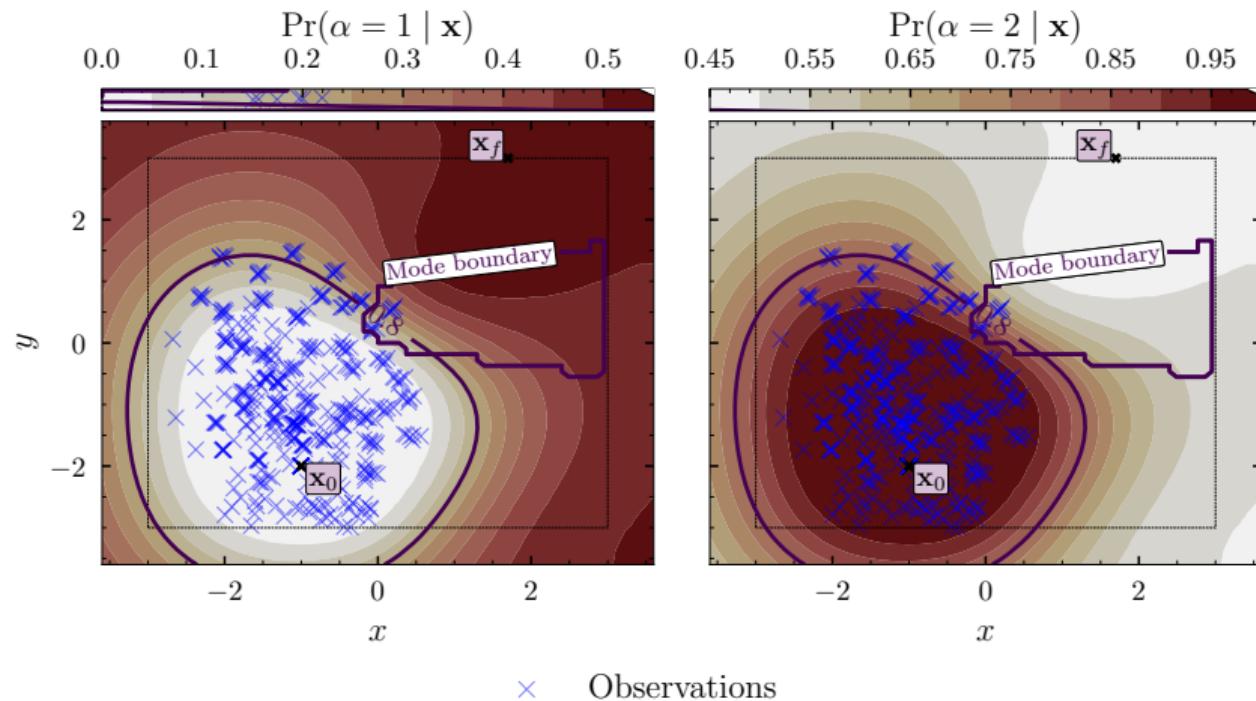
Mode remaining exploration for MBRL - iteration 3



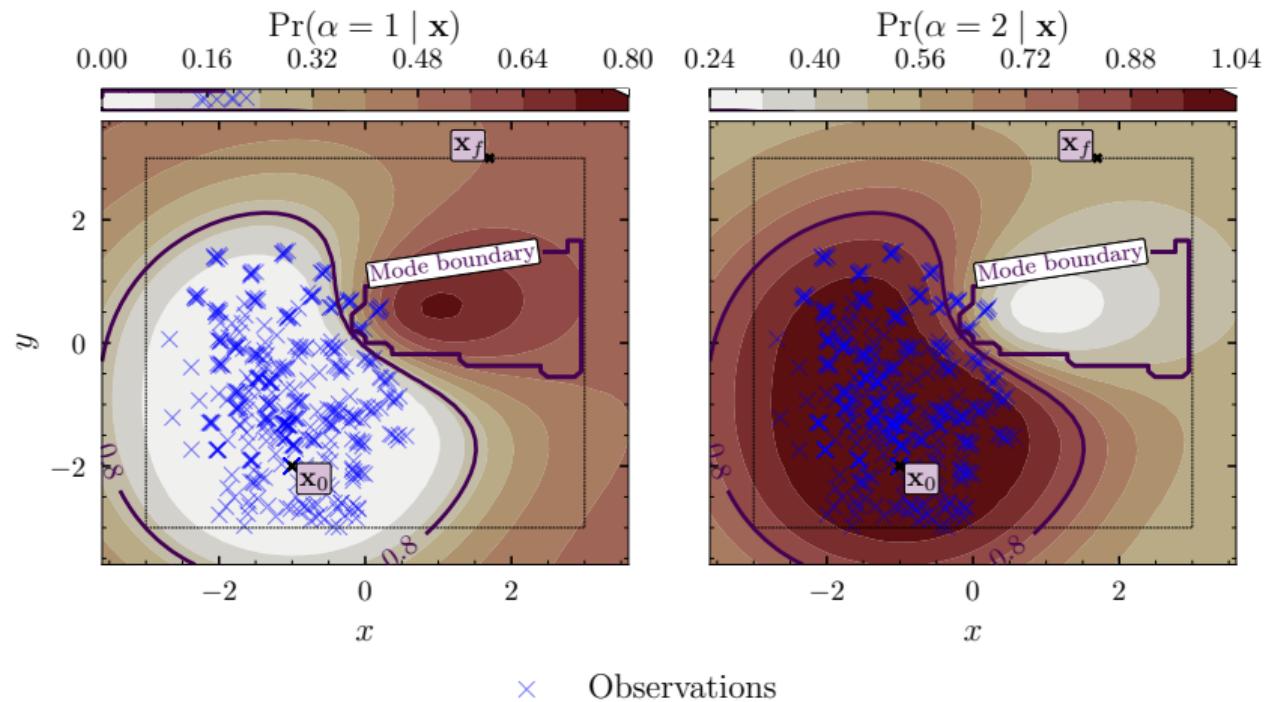
Mode remaining exploration for MBRL - iteration 3



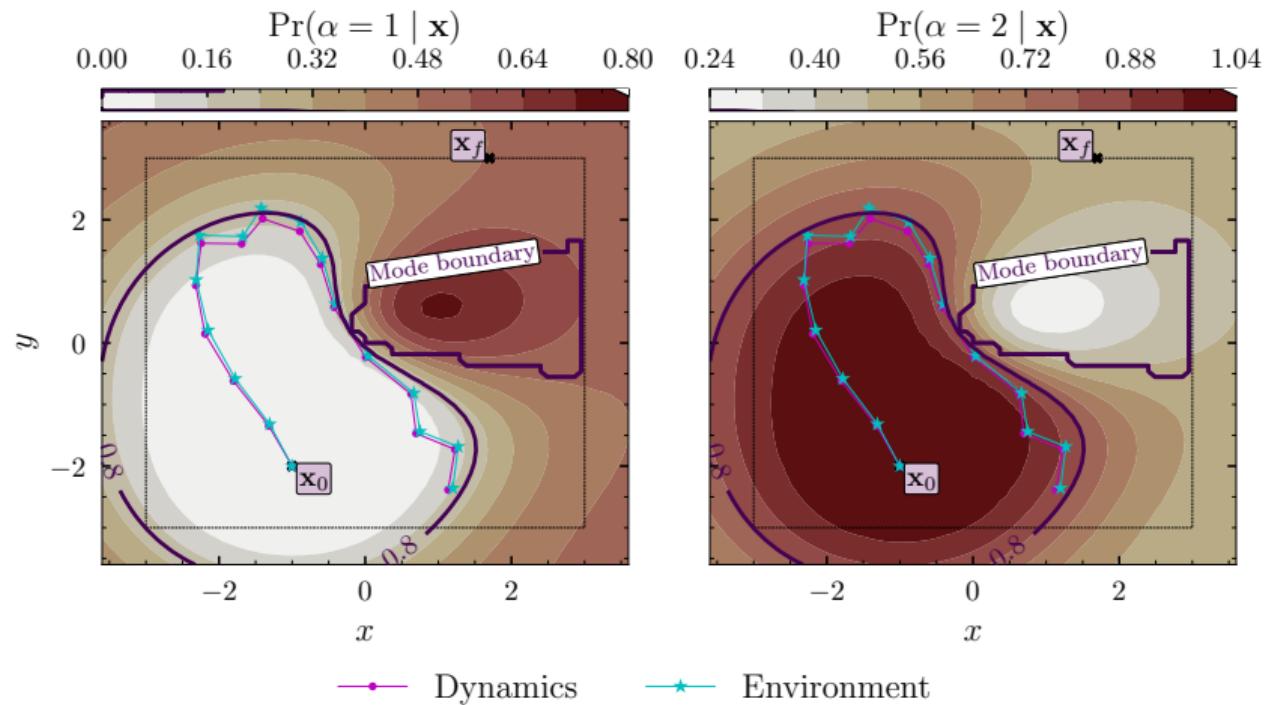
Mode remaining exploration for MBRL - iteration 3



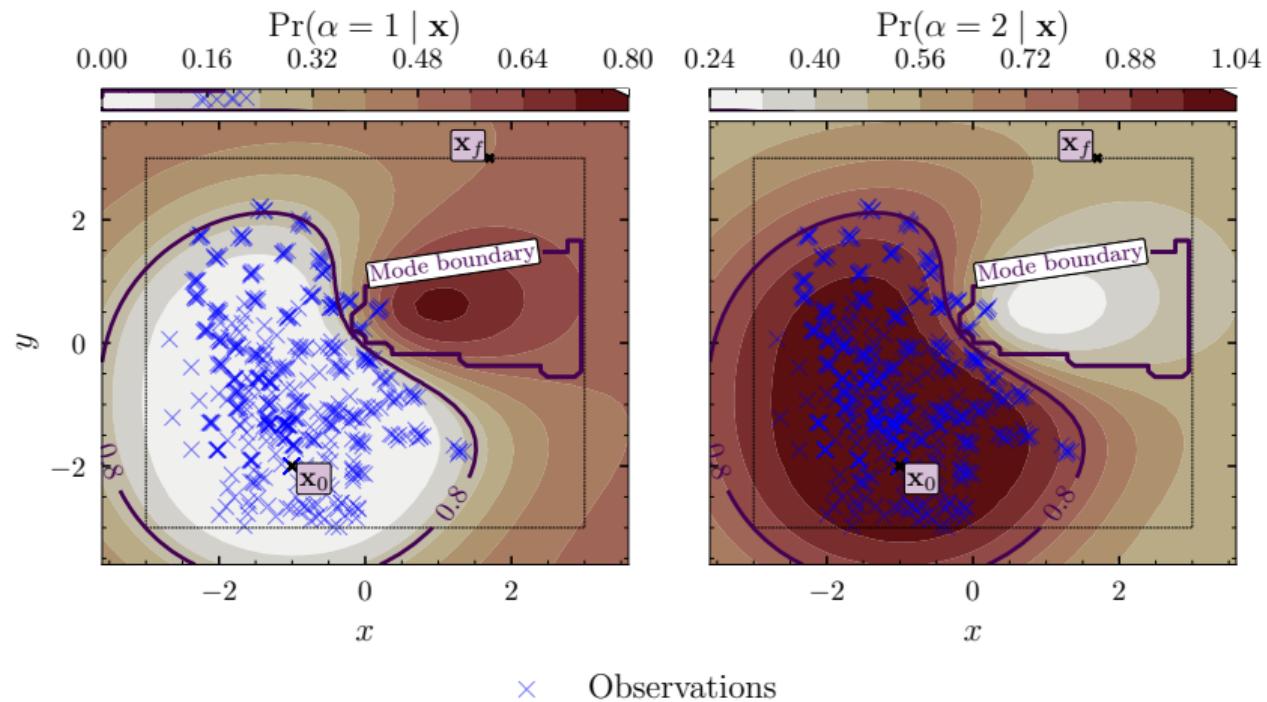
Mode remaining exploration for MBRL - iteration 4



Mode remaining exploration for MBRL - iteration 4

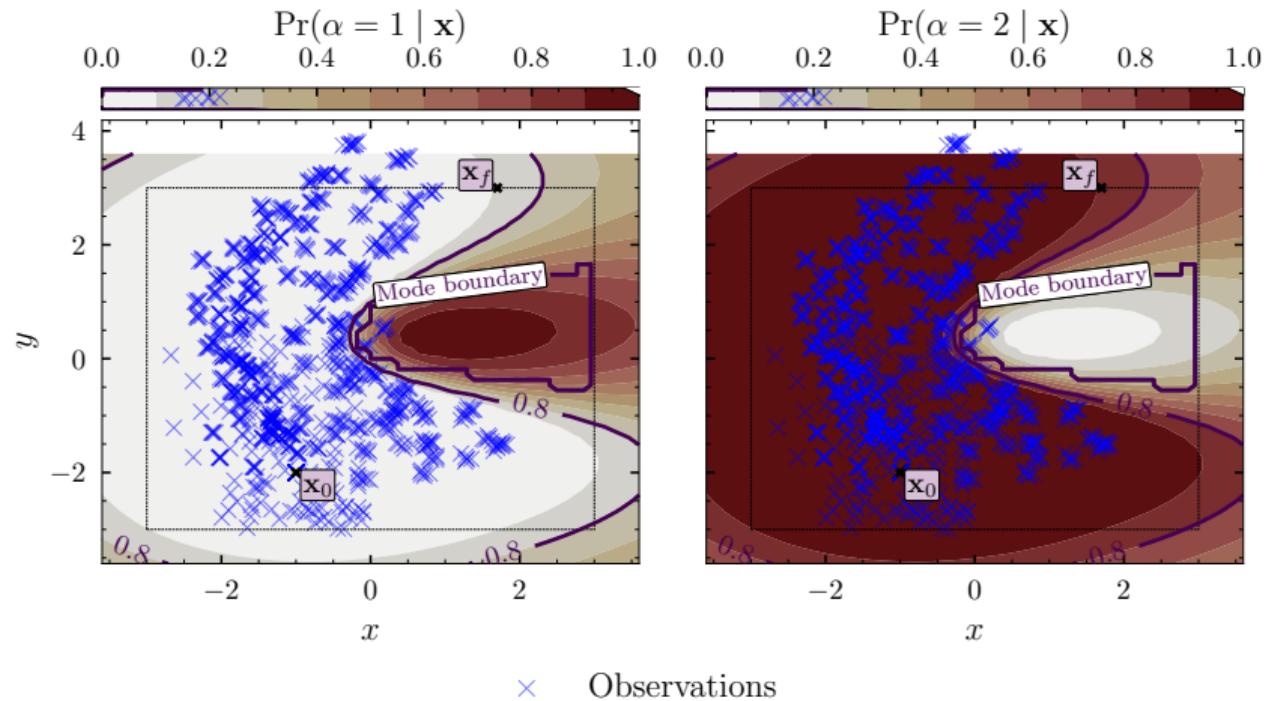


Mode remaining exploration for MBRL - iteration 4



And so on, until

And so on, until



Future work

- ❖ Bayesian treatment of inducing inputs

Future work

- ❖ Bayesian treatment of inducing inputs
- ❖ Dynamically add inducing points during exploration

Future work

- ❖ Bayesian treatment of inducing inputs
- ❖ Dynamically add inducing points during exploration
- ❖ Exploration guarantees

Future work

- ❖ Bayesian treatment of inducing inputs
- ❖ Dynamically add inducing points during exploration
- ❖ Exploration guarantees
- ❖ External sensing / higher dimensional inputs

Future work

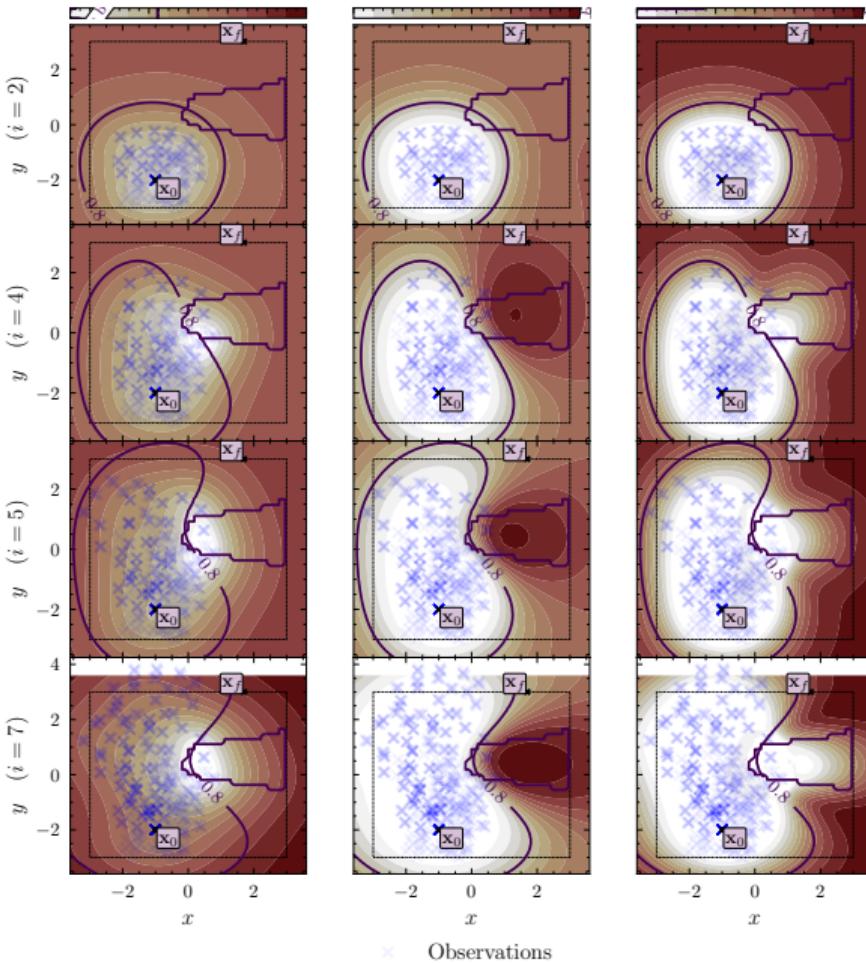
- ❖ Bayesian treatment of inducing inputs
- ❖ Dynamically add inducing points during exploration
- ❖ Exploration guarantees
- ❖ External sensing / higher dimensional inputs
- ❖ Real-time feedback control, e.g. learn a policy

Future work

- ❖ Bayesian treatment of inducing inputs
- ❖ Dynamically add inducing points during exploration
- ❖ Exploration guarantees
- ❖ External sensing / higher dimensional inputs
- ❖ Real-time feedback control, e.g. learn a policy
- ❖ Better information criterion?

$$H[h(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}_{0:i}] \quad H[\alpha \mid \mathbf{x}, \mathcal{D}_{0:i}] \quad H[\alpha \mid \mathbf{x}, \mathcal{D}_{0:i}] - \mathbb{E}_{p(h(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}_{0:i})}[H[\alpha \mid h(\mathbf{x})]]$$

0.00 1.33 2.67 4.00 5.33 6.67 8.00 0.59 0.61 0.63 0.65 0.66 0.68 0.70 0.00 0.11 0.22 0.33 0.43 0.54 0.65



What's next

- ❖ MBRL with BNN dynamics

What's next

- ❖ MBRL with BNN dynamics
 - ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?
 - ▶ Due to unimodal posterior?

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?
 - ▶ Due to unimodal posterior?
 - ▶ Idea Approximate posterior as Gaussian mixture?

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?
 - ▶ Due to unimodal posterior?
 - ▶ Idea Approximate posterior as Gaussian mixture?

❖ Multi-step dynamics models

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?
 - ▶ Due to unimodal posterior?
 - ▶ Idea Approximate posterior as Gaussian mixture?

❖ Multi-step dynamics models

- ▶ If multi-step models can outperform single-step models

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?
 - ▶ Due to unimodal posterior?
 - ▶ Idea Approximate posterior as Gaussian mixture?

❖ Multi-step dynamics models

- ▶ If multi-step models can outperform single-step models
- ▶ And, BNNs (e.g. Laplace) can work for MBRL

What's next

❖ MBRL with BNN dynamics

- ▶ Compare Laplace / MC dropout / ensembles / BNN posterior via sampling
 - ▶ Visualise *epistemic uncertainty* e.g. position vs angle for cartpole
- ▶ When do ensembles fail?
- ▶ When does Laplace approx fail?
 - ▶ Due to unimodal posterior?
 - ▶ Idea Approximate posterior as Gaussian mixture?

❖ Multi-step dynamics models

- ▶ If multi-step models can outperform single-step models
- ▶ And, BNNs (e.g. Laplace) can work for MBRL
- ▶ Idea use marginal likelihood to set multi-step model's horizon?

What's next

- ❖ Adaptivity in MBRL, i.e. switching reward functions

What's next

- ❖ Adaptivity in MBRL, i.e. switching reward functions
 - ▶ Latent dynamics learn reward functions from latent, e.g. $r_\theta : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$

What's next

- ❖ Adaptivity in MBRL, i.e. switching reward functions
 - ▶ Latent dynamics learn reward functions from latent, e.g. $r_\theta : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
 - ▶ PlaNet/DreamerV2/MuZero fail when changing reward function

What's next

- ❖ Adaptivity in MBRL, i.e. switching reward functions
 - ▶ Latent dynamics learn reward functions from latent, e.g. $r_\theta : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
 - ▶ PlaNet/DreamerV2/MuZero fail when changing reward function
 - ▶ Replay buffers lead to interference from old task

What's next

- ❖ Adaptivity in MBRL, i.e. switching reward functions
 - ▶ Latent dynamics learn reward functions from latent, e.g. $r_\theta : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
 - ▶ PlaNet/DreamerV2/MuZero fail when changing reward function
 - ▶ Replay buffers lead to interference from old task
 - ▶ Clearing replay buffer leads to catastrophic forgetting

What's next

- ☛ Adaptivity in MBRL, i.e. switching reward functions
 - ▶ Latent dynamics learn reward functions from latent, e.g. $r_\theta : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
 - ▶ PlaNet/DreamerV2/MuZero fail when changing reward function
 - ▶ Replay buffers lead to interference from old task
 - ▶ Clearing replay buffer leads to catastrophic forgetting
 - ▶ Idea Place K-priors² on dynamics?

[2] Khan et al. “Knowledge-Adaptation Priors”. 2021.

What's next

- ❖ Safety function $c : \mathcal{X} \times \mathcal{A} \rightarrow \{0 = \text{safe}, 1 = \text{unsafe}\}$

What's next

- ❖ Safety function $c : \mathcal{X} \times \mathcal{A} \rightarrow \{0 = \text{safe}, 1 = \text{unsafe}\}$
- ❖ Idea 1 Learn safety function using BNN, e.g. $c_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \{0, 1\}$

What's next

- ❖ Safety function $c : \mathcal{X} \times \mathcal{A} \rightarrow \{0 = \text{safe}, 1 = \text{unsafe}\}$
- ❖ Idea 1 Learn safety function using BNN, e.g. $c_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \{0, 1\}$
 - ▶ Probabilistic constraints

$$\prod_{t=1}^T \Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t) \geq 1 - \delta \quad (14)$$

What's next

- ❖ Safety function $c : \mathcal{X} \times \mathcal{A} \rightarrow \{0 = \text{safe}, 1 = \text{unsafe}\}$
- ❖ Idea 1 Learn safety function using BNN, e.g. $c_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \{0, 1\}$
 - ▶ Probabilistic constraints

$$\prod_{t=1}^T \Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t) \geq 1 - \delta \quad (14)$$

- ▶ which consider *epistemic uncertainty* of learned constraints function: $p(\theta \mid \mathcal{D})$

$$\Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t) = \int \Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t, \theta) p(\theta \mid \mathcal{D}) d\theta \quad (15)$$

What's next

- ❖ Safety function $c : \mathcal{X} \times \mathcal{A} \rightarrow \{0 = \text{safe}, 1 = \text{unsafe}\}$
- ❖ Idea 1 Learn safety function using BNN, e.g. $c_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \{0, 1\}$
 - ▶ Probabilistic constraints

$$\prod_{t=1}^T \Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t) \geq 1 - \delta \quad (14)$$

- ▶ which consider *epistemic uncertainty* of learned constraints function: $p(\theta \mid \mathcal{D})$

$$\Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t) = \int \Pr(c_t = 0 \mid \mathbf{x}_t, \mathbf{u}_t, \theta) p(\theta \mid \mathcal{D}) d\theta \quad (15)$$

- ▶ Can be extended to latent space dynamics models, e.g. $c_\theta : \mathcal{Z} \times \mathcal{A} \rightarrow \{0, 1\}$

What's Next

- ❖ Extending Bayesian Active Learning Disagreement³ to RL

[1] Houlsby et al. "Bayesian Active Learning for Classification and Preference Learning". 2011.

What's Next

- ❖ Extending Bayesian Active Learning Disagreement³ to RL
- ❖ In action value function space?

[1] Houlsby et al. "Bayesian Active Learning for Classification and Preference Learning". 2011.

What's Next

- ❖ Extending Bayesian Active Learning Disagreement³ to RL
- ❖ In action value function space?
 - ▶ Learn Q function $q_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

$$\pi = \arg \max_{\mathbf{u}_t} \underbrace{q_\theta(\mathbf{x}_t, \mathbf{u}_t)}_{\text{greedy}} + \underbrace{H[q_\theta(\mathbf{x}_t, \mathbf{u}_t) \mid \mathbf{x}_t, \mathbf{u}_t, \mathcal{D}_{0:i}]}_{\text{exploration}} - \mathbb{E}_{\theta \sim p(\theta \mid \mathcal{D}_{0:i})} [H[q_\theta(\mathbf{x}_t, \mathbf{u}_t) \mid \mathbf{x}_t, \mathbf{u}_t, \theta]] \quad (16)$$

[1] Houlsby et al. "Bayesian Active Learning for Classification and Preference Learning". 2011.

What's Next

- ❖ Extending Bayesian Active Learning Disagreement³ to RL
- ❖ In action value function space?
 - ▶ Learn Q function $q_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

$$\pi = \arg \max_{\mathbf{u}_t} \underbrace{q_\theta(\mathbf{x}_t, \mathbf{u}_t)}_{\text{greedy}} + \underbrace{H[q_\theta(\mathbf{x}_t, \mathbf{u}_t) \mid \mathbf{x}_t, \mathbf{u}_t, \mathcal{D}_{0:i}]}_{\text{exploration}} - \mathbb{E}_{\theta \sim p(\theta \mid \mathcal{D}_{0:i})} [H[q_\theta(\mathbf{x}_t, \mathbf{u}_t) \mid \mathbf{x}_t, \mathbf{u}_t, \theta]] \quad (16)$$

- ❖ Or in reward space?

[1] Houlsby et al. "Bayesian Active Learning for Classification and Preference Learning". 2011.

What's Next

- ❖ Extending Bayesian Active Learning Disagreement³ to RL
- ❖ In action value function space?
 - ▶ Learn Q function $q_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

$$\pi = \arg \max_{\mathbf{u}_t} \underbrace{q_\theta(\mathbf{x}_t, \mathbf{u}_t)}_{\text{greedy}} + \underbrace{H[q_\theta(\mathbf{x}_t, \mathbf{u}_t) | \mathbf{x}_t, \mathbf{u}_t, \mathcal{D}_{0:i}] - \mathbb{E}_{\theta \sim p(\theta | \mathcal{D}_{0:i})} [H[q_\theta(\mathbf{x}_t, \mathbf{u}_t) | \mathbf{x}_t, \mathbf{u}_t, \theta]]]}_{\text{exploration}}$$

(16)

- ❖ Or in reward space?

- ▶ Learn reward $r_\theta : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

$$r'(\mathbf{x}_t, \mathbf{u}_t) = \underbrace{r_\theta(\mathbf{x}_t, \mathbf{u}_t)}_{\text{greedy}} + \underbrace{H[r_\theta(\mathbf{x}_t, \mathbf{u}_t) | \mathbf{x}_t, \mathbf{u}_t, \mathcal{D}_{0:i}] - \mathbb{E}_{\theta \sim p(\theta | \mathcal{D}_{0:i})} [H[r_\theta(\mathbf{x}_t, \mathbf{u}_t) | \mathbf{x}_t, \mathbf{u}_t, \theta]]]}_{\text{exploration}}$$

(17)

[1] Houlsby et al. "Bayesian Active Learning for Classification and Preference Learning". 2011.

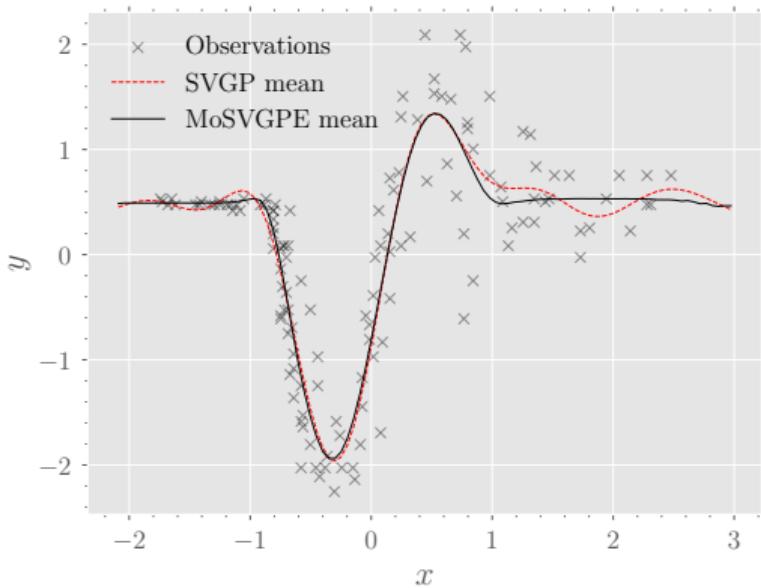
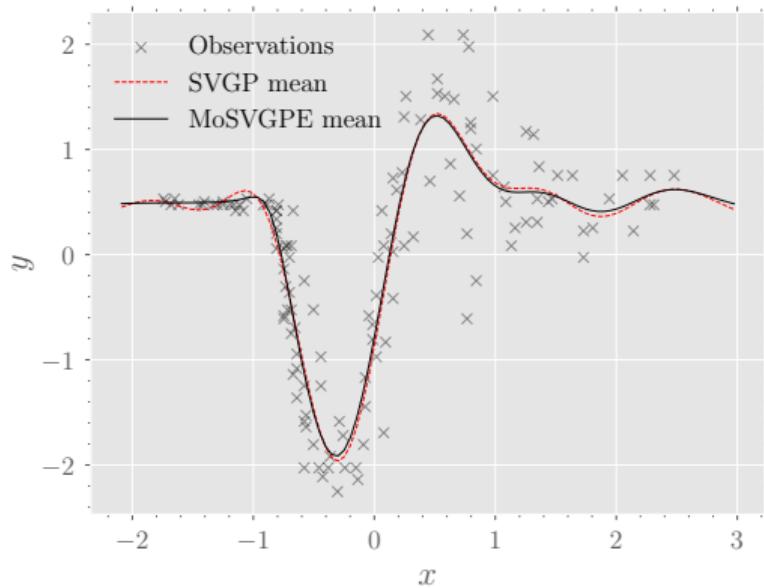
Thanks for listening

Questions?

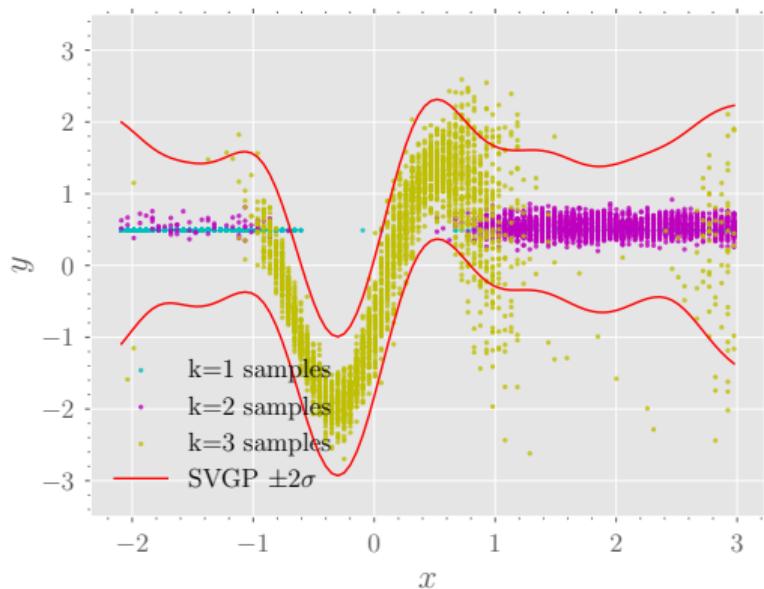
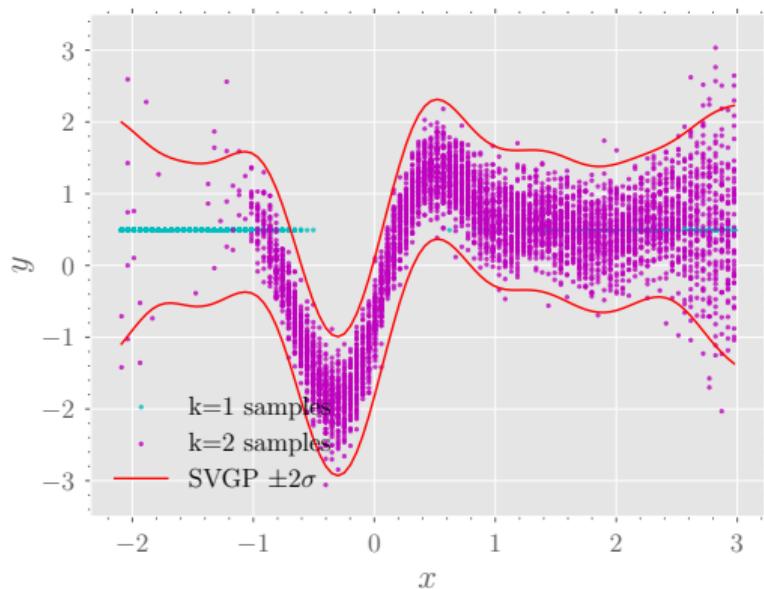
$\langle_k \text{all} \rangle$ $\langle_k \text{all} \rangle$

- [1] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. “Bayesian Active Learning for Classification and Preference Learning”. Dec. 24, 2011. arXiv: 1112.5745 [cs, stat].
- [2] Mohammad Emtiyaz E Khan and Siddharth Swaroop. “Knowledge-Adaptation Priors”. In: *Advances in Neural Information Processing Systems*. Vol. 34. Curran Associates, Inc., 2021, pp. 19757–19770.
- [3] Alessandra Tosi, Søren Hauberg, Alfredo Vellido, and Neil D Lawrence. “Metrics for Probabilistic Geometries”. In: *Proceedings of the 30th Conference*. Uncertainty in Artificial Intelligence. 2014, pp. 800–808.

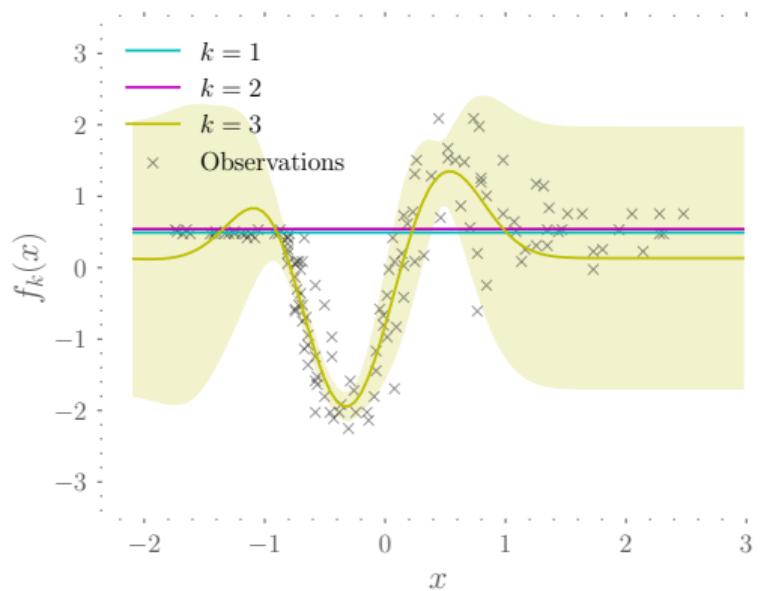
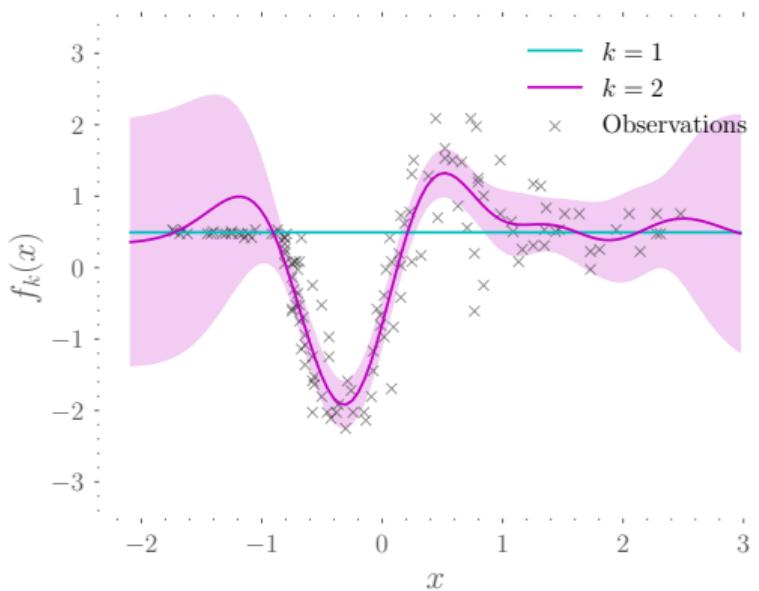
Motorcycle dataset | 2 experts vs 3 experts | Mean



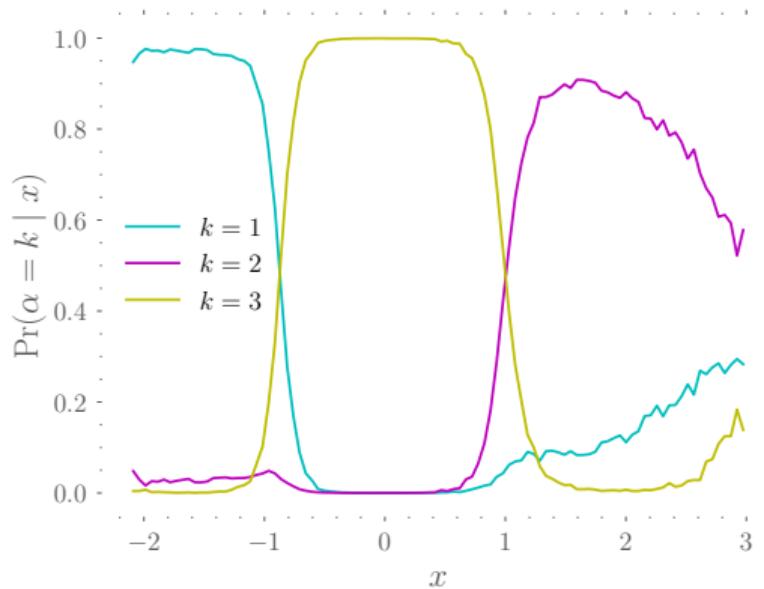
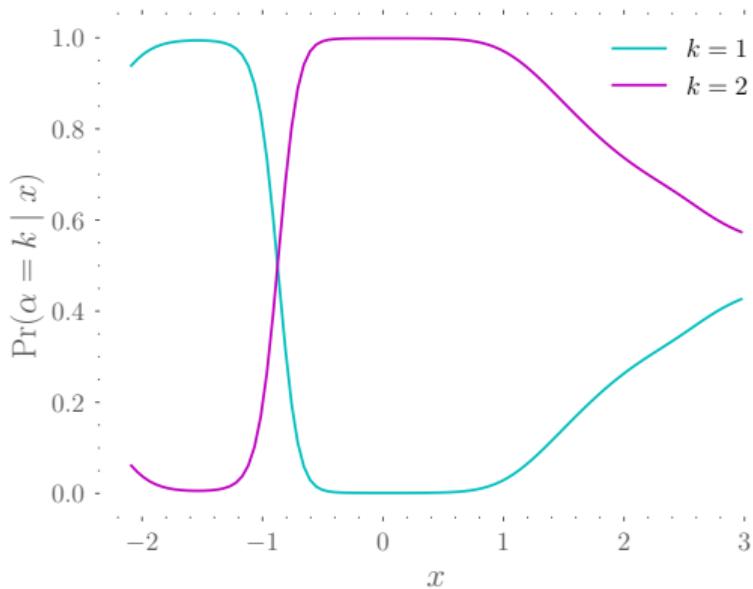
Motorcycle dataset | 2 vs 3 experts | Posterior samples



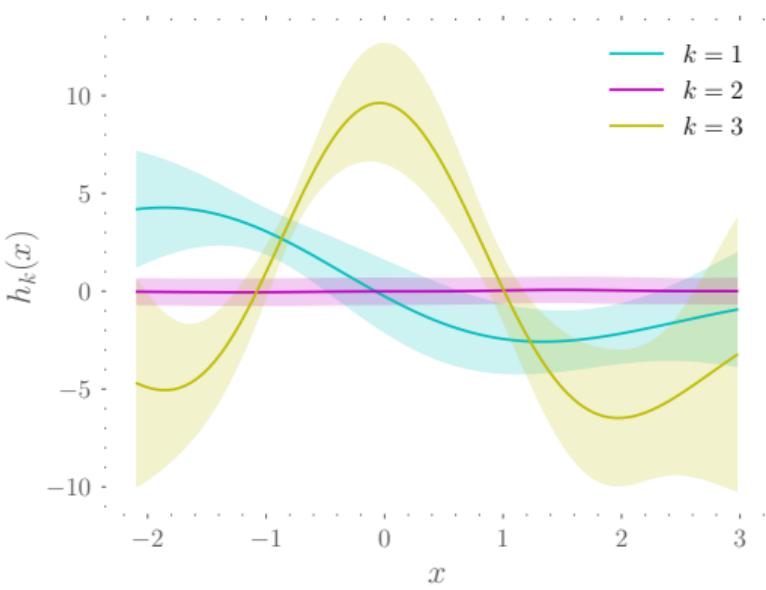
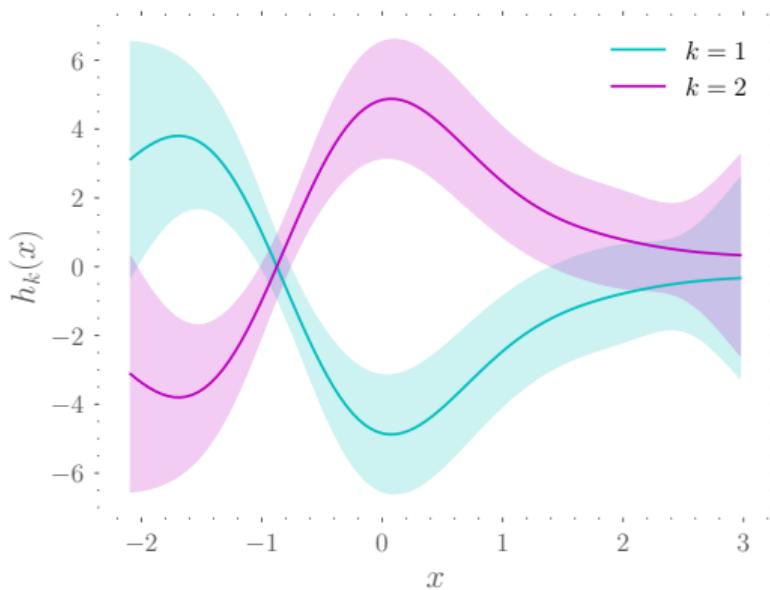
Motorcycle dataset | 2 vs 3 experts | Experts' GP posteriors



Motorcycle dataset | 2 vs 3 experts | Mixing probabilities



Motorcycle dataset | 2 vs 3 experts | Gating GP posteriors



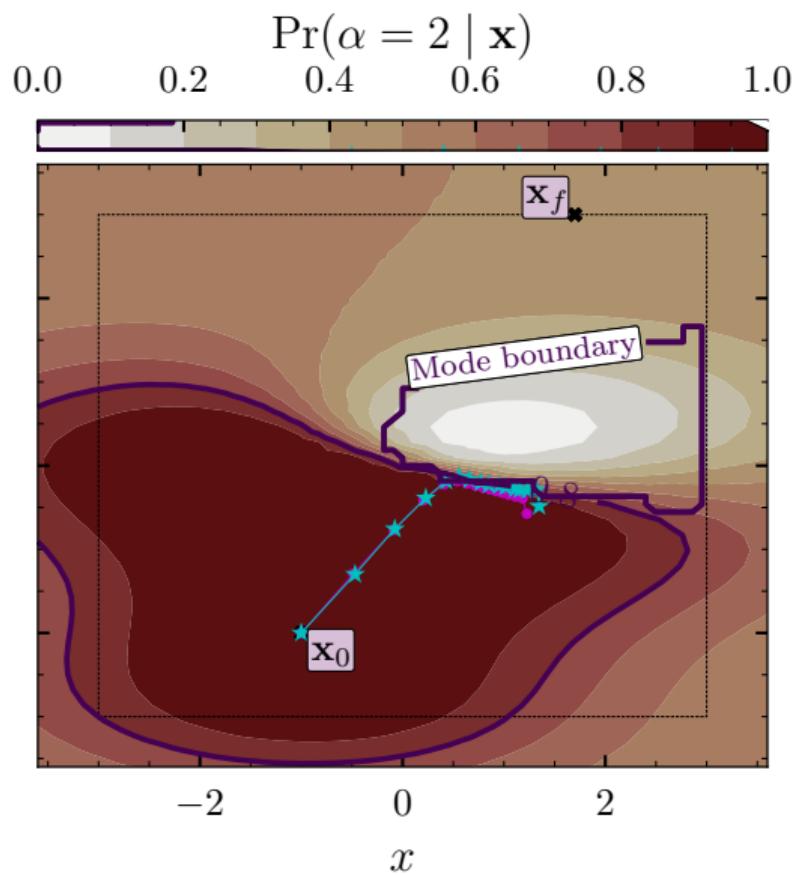


Figure: State difference only $\sum_{t=1}^{T-1} \mathbb{E} [-(\mathbf{x}_t - \mathbf{x}_f)^T \mathbf{Q} (\mathbf{x}_t - \mathbf{x}_f)]$

