

Gaussian Process Temporal Difference Learning - Theory and Practice

Yaakov Engel

Collaborators: Shie Mannor, Ron Meir, Peter Szabo,
Dmitry Volkinshtein, Nadav Aharony, Tzachi Zehavi



UNIVERSITY OF
ALBERTA

TIMELINE

- **ICML'03:** Bayes meets Bellman paper – GPTD model for MDPs with deterministic transitions
- **ICML'05:** RL with GPs paper – GPTD model for general MDPs + GPSARSA for learning control
- **NIPS'05:** Learning to control an Octopus Arm – GPTD applied to a high dimensional control problem
- **OPNET'05:** Network association-control with GPSARSA

WHY USE GPs IN RL?

- A Bayesian approach to value estimation
- Forces us to make our assumptions explicit
- Non-parametric – priors are placed and inference is performed directly in function space (kernels).
- But, can also be defined parametrically
- Domain knowledge intuitively coded in priors
- Provides full posterior, not just point estimates
- Efficient, on-line implementations, suitable for large problems

THE BAYESIAN APPROACH



- Z – hidden process, Y – observable
- We want to infer Z from measurements of Y
- Statistical dependence between Z and Y known: $P(Y|Z)$
- Place prior over Z , reflecting our uncertainty: $P(Z)$
- Observe $Y = y$
- Compute posterior: $P(Z|Y = y) = \frac{P(y|Z)P(Z)}{\int dZ' P(y|Z')P(Z')}$

GAUSSIAN PROCESSES

Definition: “An **indexed** set of jointly Gaussian random variables”

Note: The index set \mathcal{X} may be just about **any** set.

Example: $F(\mathbf{x})$, index is $\mathbf{x} \in [0, 1]^n$

F 's distribution is specified by its mean and covariance:

$$\mathbf{E}[F(\mathbf{x})] = m(\mathbf{x}) , \quad \mathbf{Cov}[F(\mathbf{x}), F(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}')$$

m is a function $\mathcal{X} \rightarrow \mathbb{R}$, k is a function $\mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$.

Conditions on k :

Symmetric, positive definite $\Rightarrow k$ is a **Mercer kernel**

GP REGRESSION

Model equation:

$$Y(\mathbf{x}) = F(\mathbf{x}) + N(\mathbf{x})$$

Prior:

$$F \sim \mathcal{N}\{0, k(\cdot, \cdot)\}$$

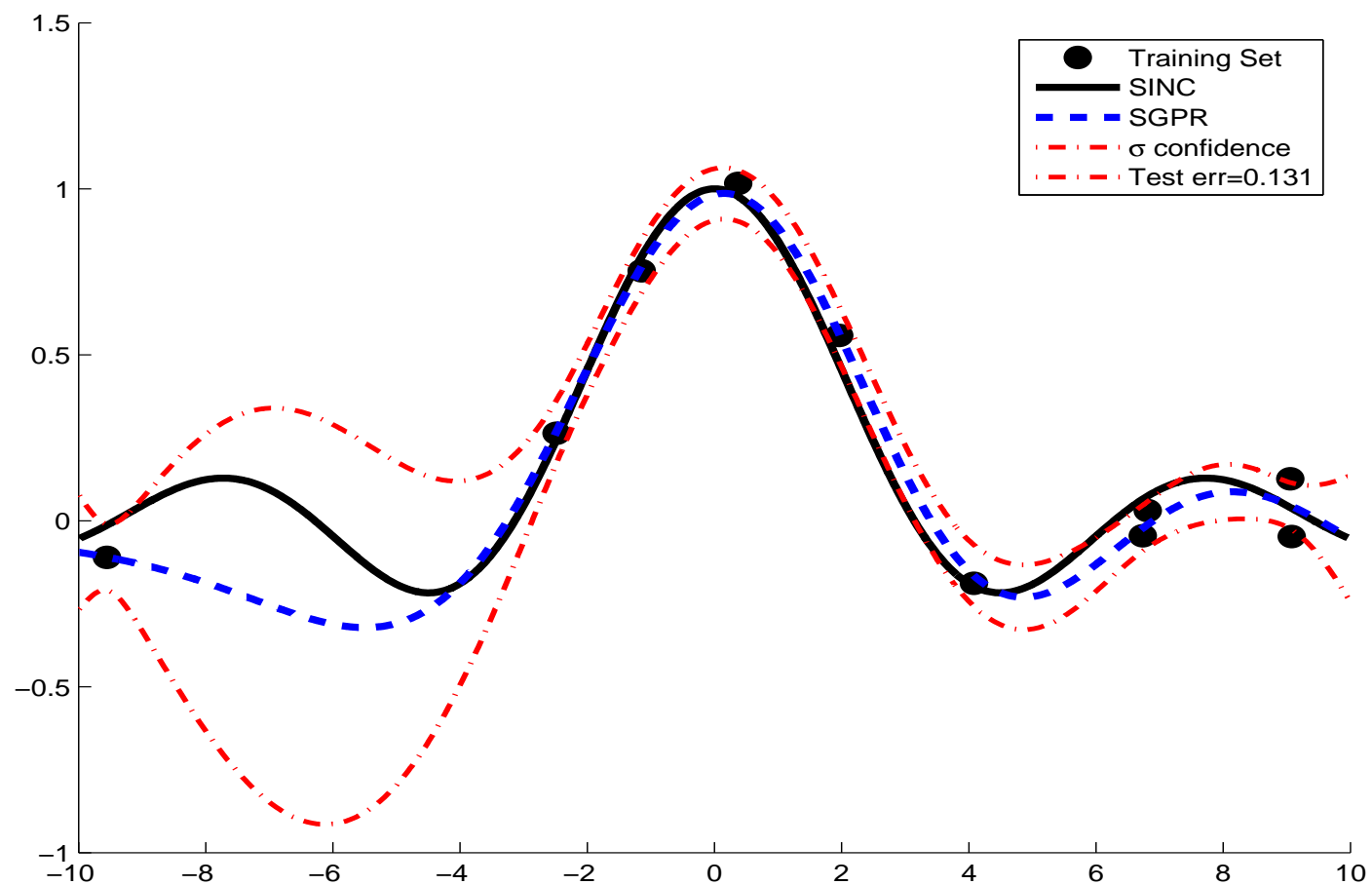
Noise:

$$N \sim \mathcal{N}\{0, \sigma^2 \delta(\cdot - \cdot)\}$$

Goal:

Find the posterior distribution of F ,
given a sample for Y (via Bayes' rule)

EXAMPLE



MARKOV DECISION PROCESSES

\mathcal{X} : state space

\mathcal{U} : action space

$p: \mathcal{X} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1], \quad \mathbf{x}_{t+1} \sim p(\cdot | \mathbf{x}_t, \mathbf{u}_t)$

$q: \mathbb{R} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1], \quad R(\mathbf{x}_t, \mathbf{u}_t) \sim q(\cdot | \mathbf{x}_t, \mathbf{u}_t)$

A Stationary policy:

$\mu: \mathcal{U} \times \mathcal{X} \rightarrow [0, 1], \quad \mathbf{u}_t \sim \mu(\cdot | \mathbf{x}_t)$

Discounted Return: $D^\mu(\mathbf{x}) = \sum_{i=0}^{\infty} \gamma^i R(\mathbf{x}_i, \mathbf{u}_i) | (\mathbf{x}_0 = \mathbf{x})$

Value function: $V^\mu(\mathbf{x}) = \mathbf{E}_\mu[D^\mu(\mathbf{x})]$

Goal: Find a policy μ^* maximizing $V^\mu(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}$

BELLMAN'S EQUATION

For a fixed policy μ :

$$V^\mu(\mathbf{x}) = \mathbf{E}_{\mathbf{x}', \mathbf{u} | \mathbf{x}} \left[R(\mathbf{x}, \mathbf{u}) + \gamma V^\mu(\mathbf{x}') \right]$$

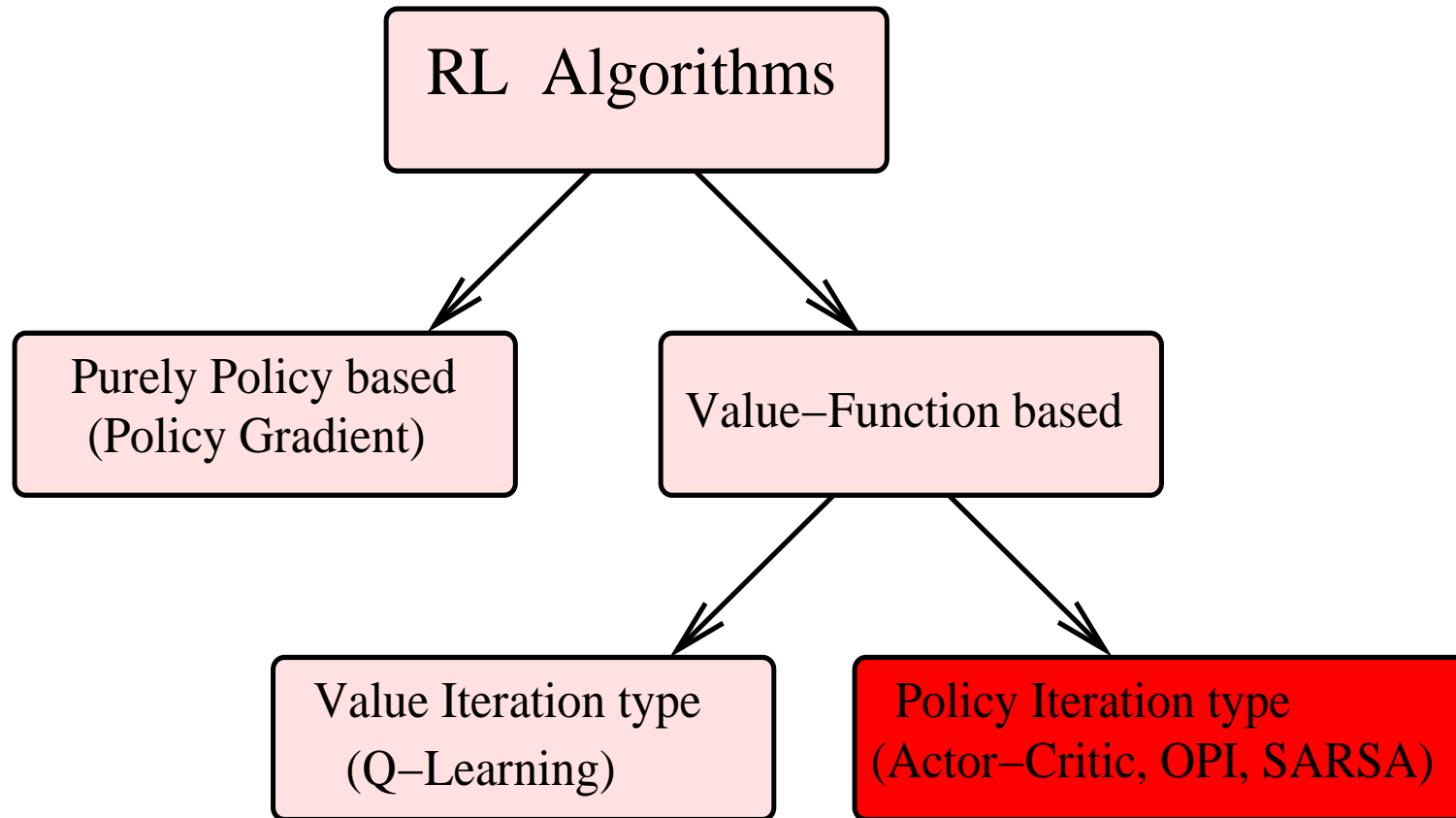
Optimal value and policy:

$$V^*(\mathbf{x}) = \max_{\mu} V^\mu(\mathbf{x}) , \quad \mu^* = \operatorname{argmax}_{\mu} V^\mu(\mathbf{x})$$

How to solve it?

- Methods based on Value Iteration (e.g. Q-learning)
- Methods based on Policy Iteration (e.g. SARSA, OPI, Actor-Critic)

SOLUTION METHOD TAXONOMY



PI methods need a “subroutine” for policy evaluation

WHAT'S MISSING?

Shortcomings of current policy evaluation methods:

- Some methods can only be applied to small problems
- No probabilistic interpretation - how good is the estimate?
- Only parametric methods are capable of operating on-line
- Non-parametric methods are more flexible but only work off-line
- Small-step-size (stoch. approx.) methods use data inefficiently
- Finite-time solutions lack interpretability, all statements are asymptotic
- Convergence issues

GAUSSIAN PROCESS TEMPORAL DIFFERENCE LEARNING

Model Equations:

$$R(\mathbf{x}_i) = V(\mathbf{x}_i) - \gamma V(\mathbf{x}_{i+1}) + N(\mathbf{x}_i, \mathbf{x}_{i+1})$$

Or, in compact form:

$$R_t = \mathbf{H}_{t+1} V_{t+1} + N_t$$

$$\mathbf{H}_t = \begin{bmatrix} 1 & -\gamma & 0 & \dots & 0 \\ 0 & 1 & -\gamma & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & \dots & 1 & -\gamma \end{bmatrix}.$$

Our (Bayesian) goal:

Find the posterior distribution of V ,
given a sequence of observed states and rewards.

DETERMINISTIC DYNAMICS

Bellman's Equation:

$$V(\mathbf{x}_i) = \bar{R}(\mathbf{x}_i) + \gamma V(\mathbf{x}_{i+1})$$

Define:

$$N(\mathbf{x}) = R(\mathbf{x}) - \bar{R}(\mathbf{x})$$

Assumption: $N(\mathbf{x}_i)$ are Normal, i.i.d., with variance σ^2 .

Model Equations:

$$R(\mathbf{x}_i) = V(\mathbf{x}_i) - \gamma V(\mathbf{x}_{i+1}) + N(\mathbf{x}_i)$$

In compact form:

$$R_t = \mathbf{H}_{t+1} V_{t+1} + N_t, \text{ with } N_t \sim \mathcal{N}\{0, \sigma^2 \mathbf{I}\}$$

STOCHASTIC DYNAMICS

The discounted return:

$$D(\mathbf{x}_i) = \mathbf{E}_\mu D(\mathbf{x}_i) + (D(\mathbf{x}_i) - \mathbf{E}_\mu D(\mathbf{x}_i)) = V(\mathbf{x}_i) + \Delta V(\mathbf{x}_i)$$

For a stationary MDP:

$$D(\mathbf{x}_i) = R(\mathbf{x}_i) + \gamma D(\mathbf{x}_{i+1}) \text{ (where } \mathbf{x}_{i+1} \sim p(\cdot|\mathbf{x}_i, \mathbf{u}_i), \mathbf{u}_i \sim \mu(\cdot|\mathbf{x}_i))$$

Substitute and rearrange:

$$\begin{aligned} R(\mathbf{x}_i) &= V(\mathbf{x}_i) - \gamma V(\mathbf{x}_{i+1}) + N(\mathbf{x}_i, \mathbf{x}_{i+1}) \\ N(\mathbf{x}_i, \mathbf{x}_{i+1}) &\stackrel{\text{def}}{=} \Delta V(\mathbf{x}_i) - \gamma \Delta V(\mathbf{x}_{i+1}) \end{aligned}$$

Assumption: $\Delta V(\mathbf{x}_i)$ are Normal, i.i.d., with variance σ^2 .

In compact form:

$$R_t = \mathbf{H}_{t+1} V_{t+1} + N_t, \text{ with } N_t \sim \mathcal{N}\{0, \sigma^2 \mathbf{H}_{t+1} \mathbf{H}_{t+1}^\top\}$$

THE POSTERIOR

General noise covariance:

$$\text{Cov}[N_t] = \Sigma_t$$

Joint distribution:

$$\begin{bmatrix} R_{t-1} \\ V(\mathbf{x}) \end{bmatrix} \sim \mathcal{N} \left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{H}_t \mathbf{K}_t \mathbf{H}_t^\top + \Sigma_t & \mathbf{H}_t \mathbf{k}_t(\mathbf{x}) \\ \mathbf{k}_t(\mathbf{x})^\top \mathbf{H}_t^\top & k(\mathbf{x}, \mathbf{x}) \end{bmatrix} \right\}$$

Invoke Bayes' Rule:

$$\mathbf{E}[V(\mathbf{x}) | R_{t-1} = \mathbf{r}_{t-1}] = \mathbf{k}_t(\mathbf{x})^\top \boldsymbol{\alpha}_t$$

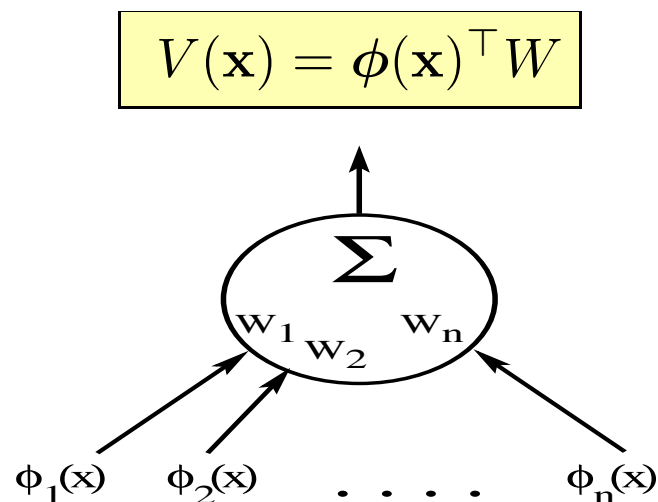
$$\text{Cov}[V(\mathbf{x}), V(\mathbf{x}') | R_{t-1} = \mathbf{r}_{t-1}] = k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_t(\mathbf{x})^\top \mathbf{C}_t \mathbf{k}_t(\mathbf{x}')$$

$$\mathbf{k}_t(\mathbf{x}) = (k(\mathbf{x}_0, \mathbf{x}), \dots, k(\mathbf{x}_t, \mathbf{x}))^\top, \quad \mathbf{K}_t = [\mathbf{k}_t(\mathbf{x}_0), \dots, \mathbf{k}_t(\mathbf{x}_t)]$$

$$\boldsymbol{\alpha}_t = \mathbf{H}_t^\top \left(\mathbf{H}_t \mathbf{K}_t \mathbf{H}_t^\top + \Sigma_t \right)^{-1} \mathbf{r}_{t-1}, \quad \mathbf{C}_t = \mathbf{H}_t^\top \left(\mathbf{H}_t \mathbf{K}_t \mathbf{H}_t^\top + \Sigma_t \right)^{-1} \mathbf{H}_t.$$

A PARAMETRIC GAUSSIAN PROCESS MODEL

A linear combination of features:



Prior on W : Gaussian, with $\mathbf{E}[W] = \mathbf{0}$, $\text{Cov}[W, W] = \mathbf{I}$

Prior on V : Gaussian, with

$$\mathbf{E}[V(\mathbf{x})] = \mathbf{0}, \quad \text{Cov}[V(\mathbf{x}), V(\mathbf{x}')] = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

COMPARISON OF MODELS

	Parametric	Nonparametric
Parametrization	$V(\mathbf{x}) = \phi(\mathbf{x})^\top W$	None, V is V
Prior	$W \sim \mathcal{N}\{\mathbf{0}, \mathbf{I}\}$	$V \sim \mathcal{N}\{0, k(\cdot, \cdot)\}$
$\mathbf{E}[V(\mathbf{x})]$	0	0
$\mathbf{Cov}[V(\mathbf{x}), V(\mathbf{x}')]]$	$\phi(\mathbf{x})^\top \phi(\mathbf{x}')$	$k(\mathbf{x}, \mathbf{x}')$
We seek	$W R_{t-1}$	$V(\mathbf{x}) R_{t-1}$

If we can find a set of basis functions satisfying $\phi(\mathbf{x})^\top \phi(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$, the two models become equivalent.

In fact, such a set **always** exists [Mercer].

However, it may be infinite

RELATION TO MONTE-CARLO ESTIMATION

In the stochastic model: $\Sigma_t = \sigma^2 \mathbf{H}_{t+1} \mathbf{H}_{t+1}^\top$

Also, let: $(Y_t)_i = \sum_{j=i}^t \gamma^{j-i} R(\mathbf{x}_j, \mathbf{u}_j)$

Then:

$$\begin{aligned}\mathbf{E}[W|R_t] &= \left(\Phi_t \Phi_t^\top + \sigma^2 \mathbf{I} \right)^{-1} \Phi_t Y_t \\ \text{Cov}[W|R_t] &= \sigma^2 \left(\Phi_t \Phi_t^\top + \sigma^2 \mathbf{I} \right)^{-1}\end{aligned}$$

That's the solution to GP regression on Monte-Carlo samples of the discounted return.

MAP / ML SOLUTIONS

Since the posterior is Gaussian:

$$\hat{\mathbf{w}}_{t+1}^{MAP} = \mathbf{E}[W|R_t] = \left(\Phi_t \Phi_t^\top + \sigma^2 \mathbf{I} \right)^{-1} \Phi_t Y_t$$

Performing ML inference using the same model we get:

$$\hat{\mathbf{w}}_{t+1}^{ML} = \left(\Phi_t \Phi_t^\top \right)^{-1} \Phi_t Y_t$$

That's the LSTD(1) (Least-Squares Monte-Carlo) solution.

POLICY IMPROVEMENT

How can we perform policy improvement?

State values? Not without a transition model (even then tricky).

State-action (Q-) values? Yes!

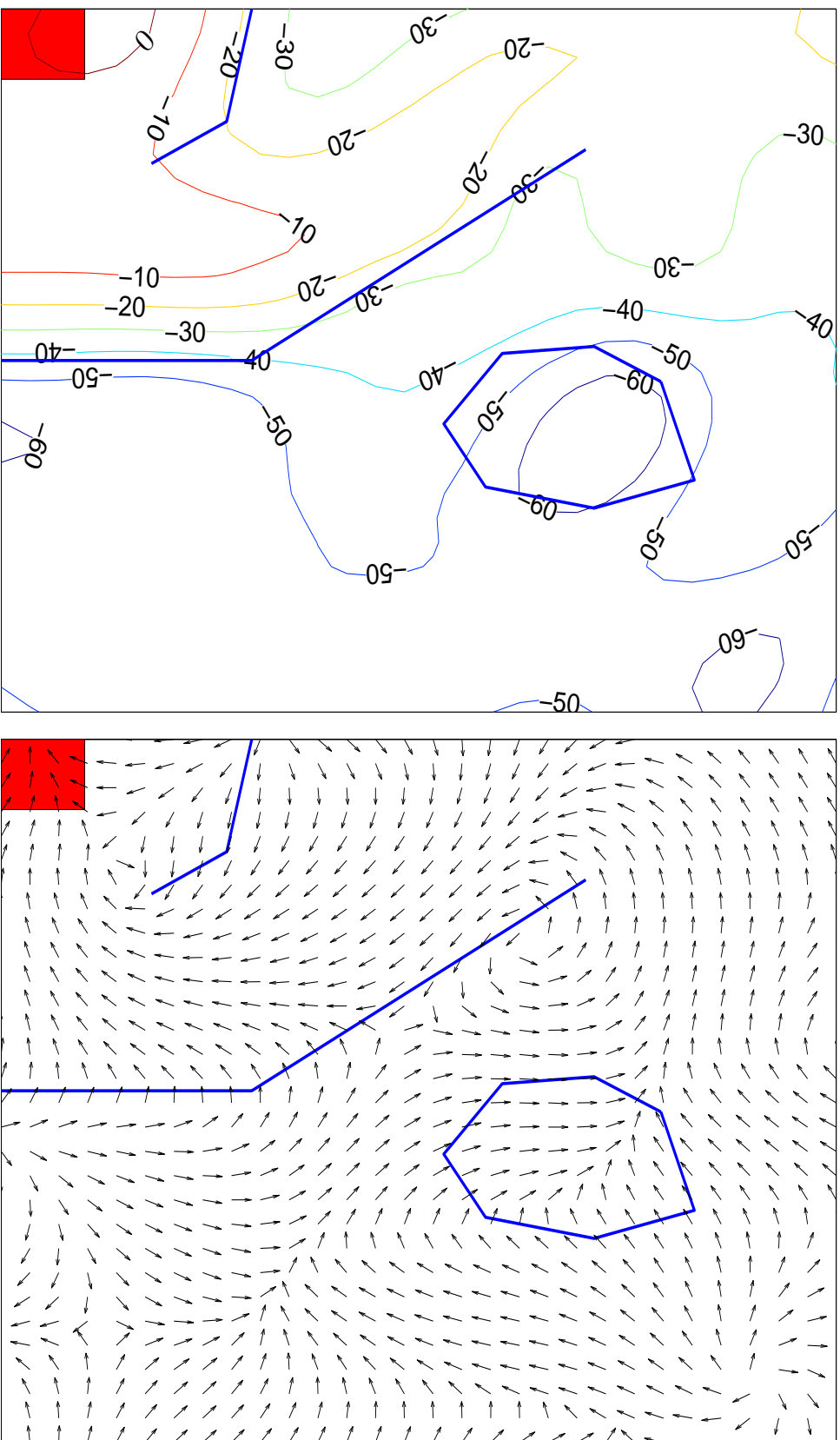
Idea: Use a state-action value GP

How?

- Define a state-action kernel: $k((\mathbf{x}, \mathbf{u}), (\mathbf{x}', \mathbf{u}'))$
- Run GPTD on state-action pairs
- Use some semi-greedy action selection rule

We call this GPSARSA.

A SIMPLE EXPERIMENT



THE OCTOPUS ARM

Can bend and twist at any point

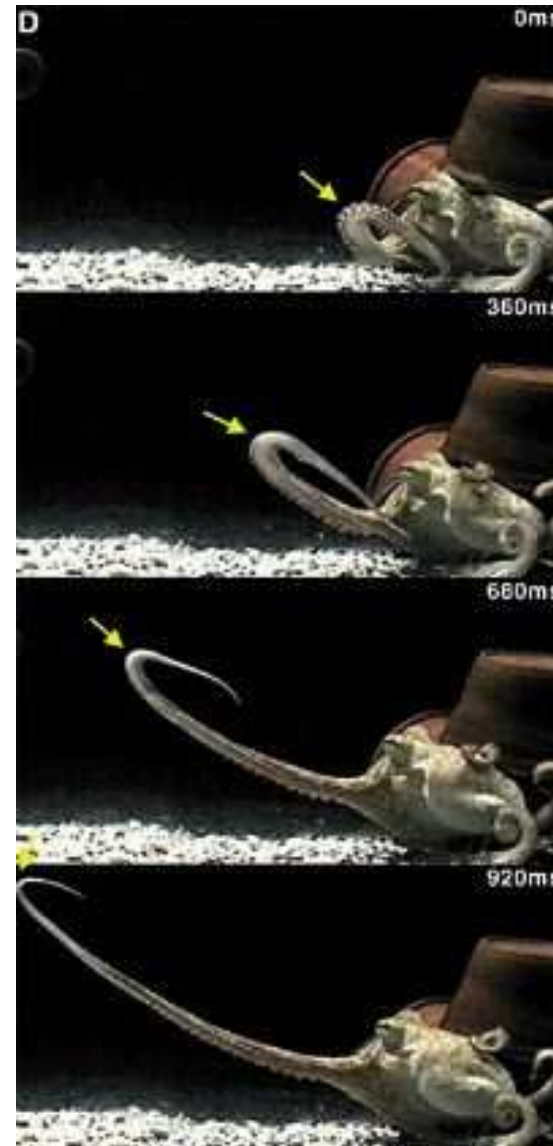
Can do this in any direction

Can be elongated and shortened

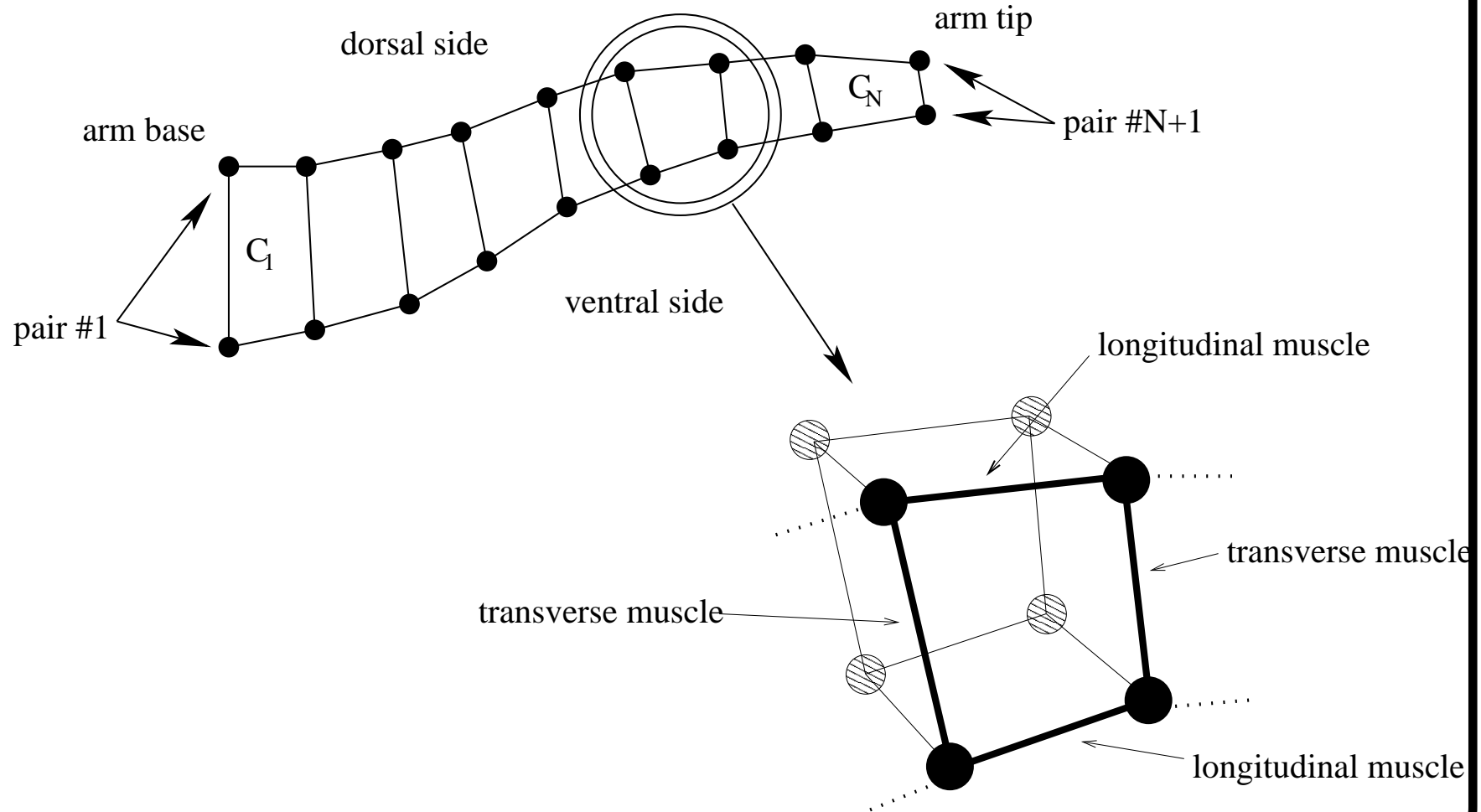
Can change cross section

Can grab using any part of the arm

Virtually infinitely many DOF

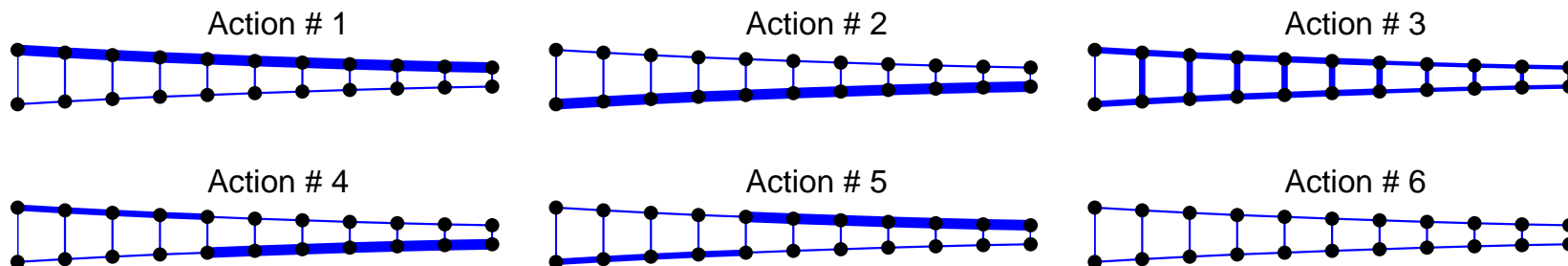


OUR ARM MODEL



ACTIONS

Each action specifies a set of fixed activations – one for each muscle in the arm.



Base rotation adds duplicates of actions 1,2,4 and 5 with positive and negative torques applied to the base.

THE CONTROL PROBLEM

Starting from a random position, bring {any part, tip} of arm into contact with a goal region, **optimally**.

Optimality criteria:

Time, energy, obstacle avoidance

Constraint:

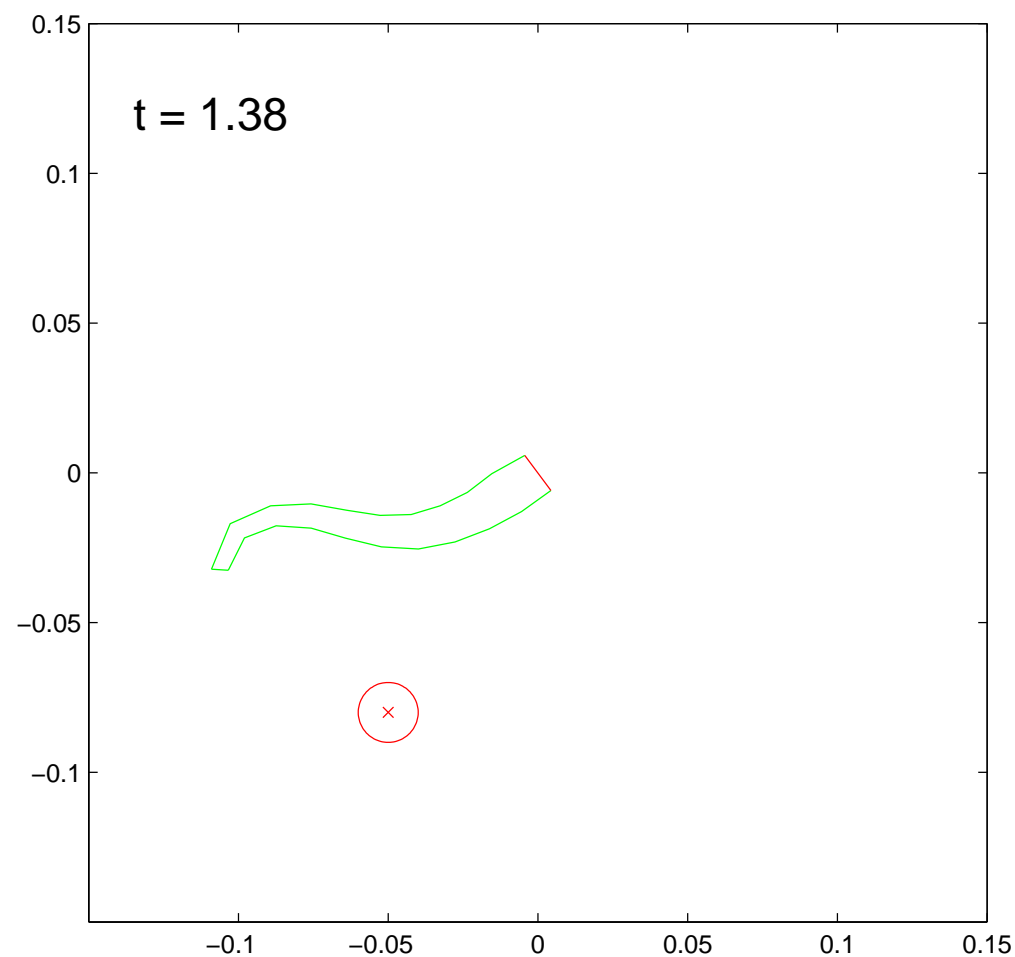
We only have access to sampled trajectories

Our approach:

Define problem as a MDP

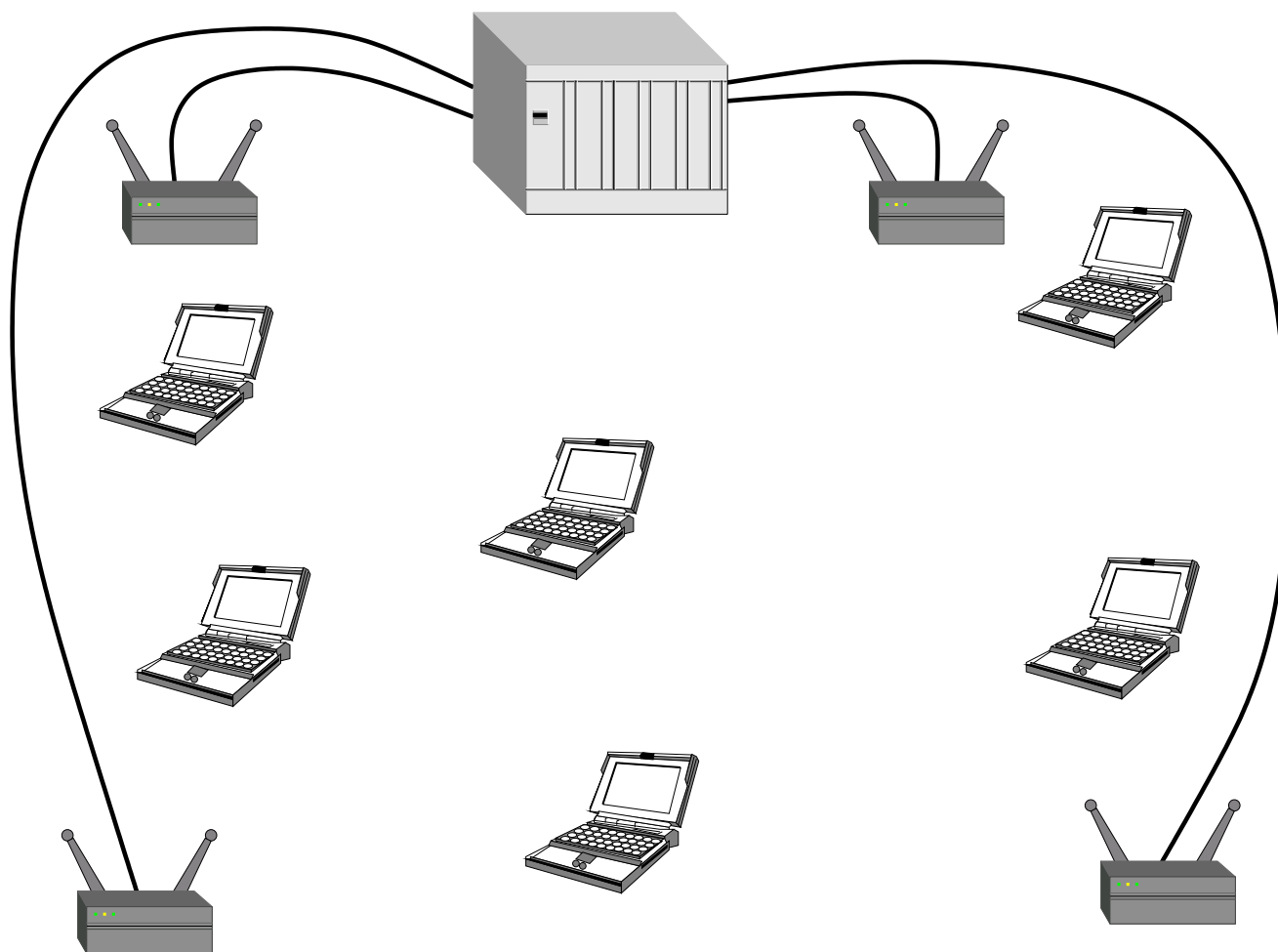
Solve using a GPTD algorithm

THE TASK



MOVIES

ASSOCIATION CONTROL IN WLANs



ASSOCIATION CONTROL IN WLANs

Setting: n users, $m \ll n$ access points (APs),

The problem: Associate users with APs, optimally.

Complications: Users are not the same, they move around, change their behavior over time, what is meant by “optimally”? etc.

Idea: Model the system as a MDP, solve using GPSARSA

Results:

- Tested on simple networks using the OPNET simulator
- Preliminary results look promising
- More work is needed

CHALLENGES

- How to use value uncertainty?
- What's a disciplined way to select actions?
- What's the best noise covariance?
- Bias, variance, learning curves
- POMDPs
- More complicated tasks

Questions?