

Bayes



RL

Gaussian Process Temporal Difference Learning

Yaakov Engel

Collaborators: Shie Mannor, Ron Meir



Bayes



RL

WHY USE GPs IN RL?

- A Bayesian approach to value estimation
- Forces us to 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 over values, not just point estimates
- Efficient, on-line implementations, suitable for large problems



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:

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

Conditions on k :

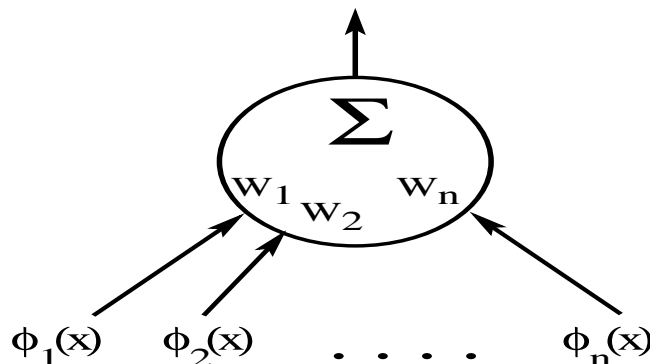
Symmetric, positive definite $\Rightarrow k$ is a **Mercer kernel**



EXAMPLE: PARAMETRIC GP

A linear combination of basis functions:

$$F(\mathbf{x}) = \phi(\mathbf{x})^\top W$$



If $W \sim \mathcal{N}\{\mathbf{m}_w, \mathbf{C}_w\}$,

then F is a GP:

$$\begin{aligned} \mathbf{E}[F(\mathbf{x})] &= \phi(\mathbf{x})^\top \mathbf{m}_w, \\ \text{Cov}[F(\mathbf{x}), F(\mathbf{x}')] &= \phi(\mathbf{x})^\top \mathbf{C}_w \phi(\mathbf{x}') \end{aligned}$$



CONDITIONING – GAUSS-MARKOV THM.

Theorem Let Z and Y be random vectors jointly distributed according to the multivariate normal distribution

$$\begin{pmatrix} Z \\ Y \end{pmatrix} \sim \mathcal{N} \left\{ \begin{pmatrix} \mathbf{m}_z \\ \mathbf{m}_y \end{pmatrix}, \begin{bmatrix} \mathbf{C}_{zz} & \mathbf{C}_{zy} \\ \mathbf{C}_{yz} & \mathbf{C}_{yy} \end{bmatrix} \right\}.$$

Then $Z|Y \sim \mathcal{N} \{ \hat{Z}, \mathbf{P} \}$, where

$$\hat{Z} = \mathbf{m}_z + \mathbf{C}_{zy} \mathbf{C}_{yy}^{-1} (Y - \mathbf{m}_y)$$

$$\mathbf{P} = \mathbf{C}_{zz} - \mathbf{C}_{zy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{yz}.$$

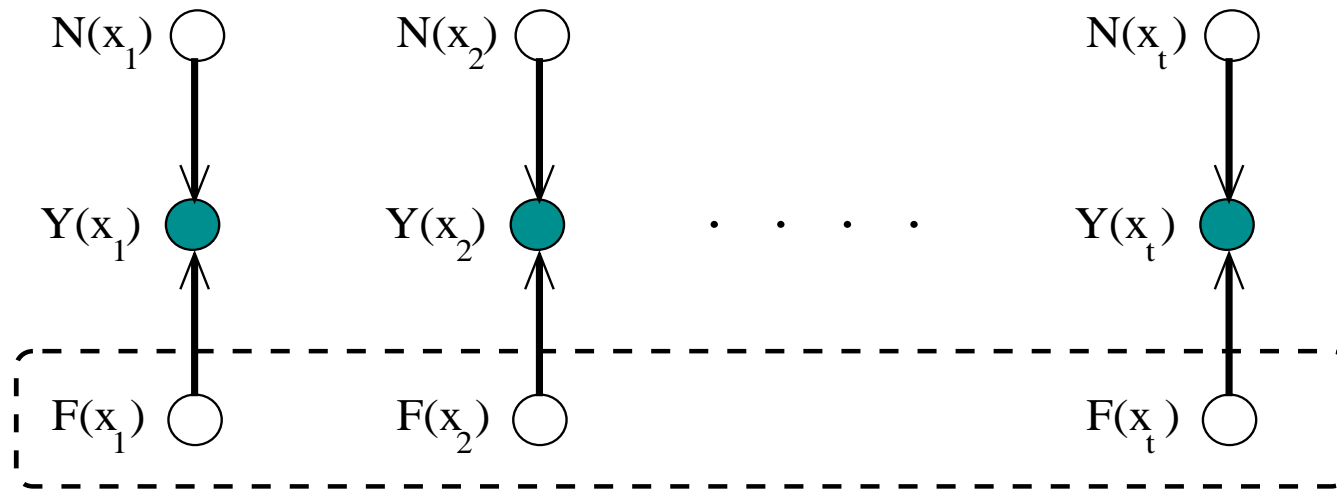


GP REGRESSION

Sample: $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_t, y_t))$

Model equation: $Y(\mathbf{x}_i) = F(\mathbf{x}_i) + N(\mathbf{x}_i)$

GP Prior on F : $F \sim \mathcal{N}\{0, k(\cdot, \cdot)\}$



N : IID zero-mean Gaussian noise, with variance σ^2



GP REGRESSION (CTD.)

Denote:

$$Y_t = (Y(\mathbf{x}_1), \dots, Y(\mathbf{x}_t))^{\top},$$

$$\mathbf{k}_t(\mathbf{x}) = (k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_t, \mathbf{x}))^{\top},$$

$$\mathbf{K}_t = [\mathbf{k}_t(\mathbf{x}_1), \dots, \mathbf{k}_t(\mathbf{x}_t)].$$

Then:

$$\begin{pmatrix} F(\mathbf{x}) \\ Y_t \end{pmatrix} \sim \mathcal{N} \left\{ \begin{pmatrix} 0 \\ \mathbf{0} \end{pmatrix}, \begin{bmatrix} k(\mathbf{x}, \mathbf{x}) & \mathbf{k}_t(\mathbf{x})^{\top} \\ \mathbf{k}_t(\mathbf{x}) & \mathbf{K}_t + \sigma^2 \mathbf{I} \end{bmatrix} \right\}$$

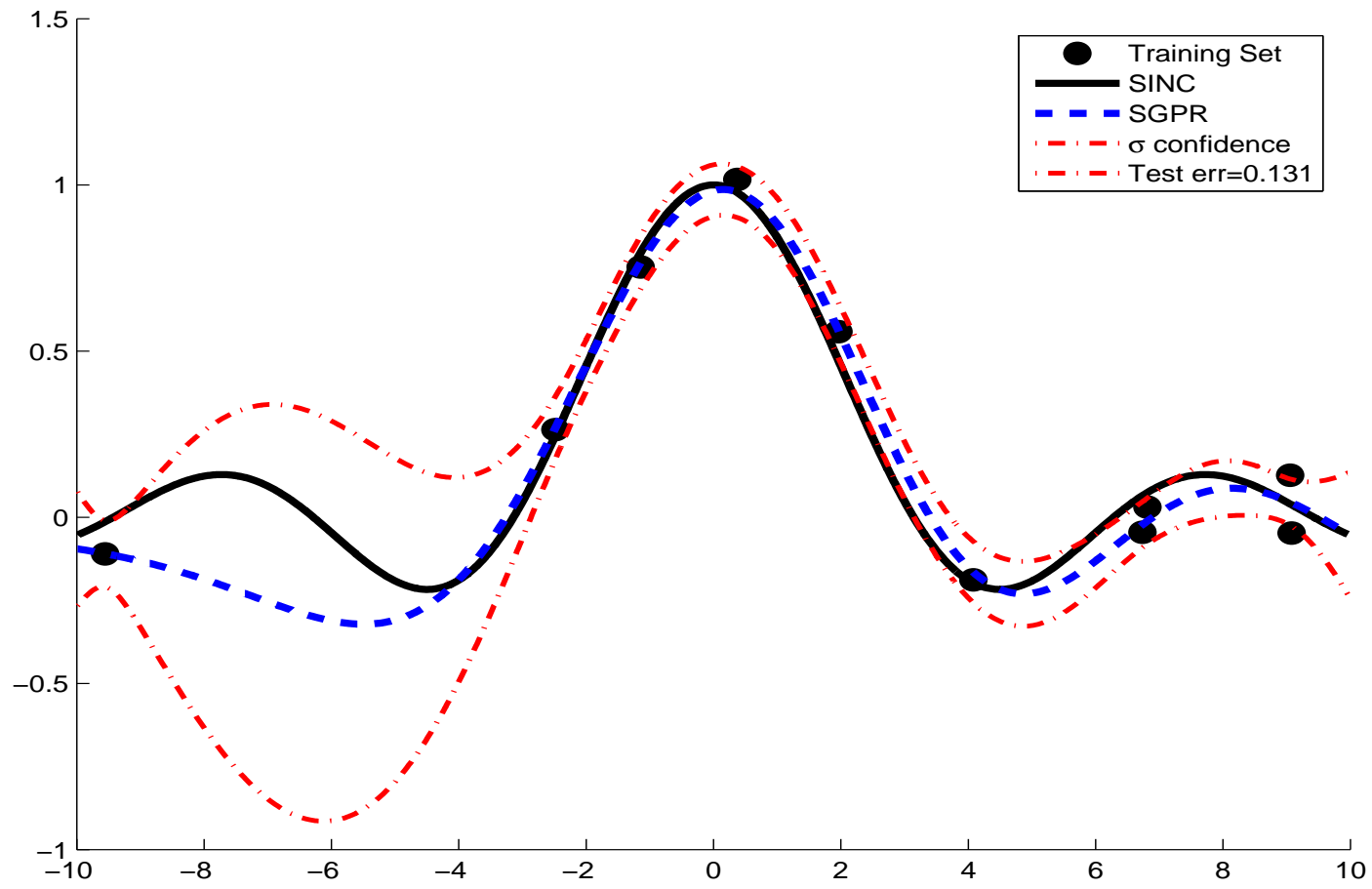
Now apply conditioning formula to compute the posterior moments of $F(\mathbf{x})$, given $Y_t = \mathbf{y}_t = (y_1, \dots, y_t)^{\top}$.

Bayes



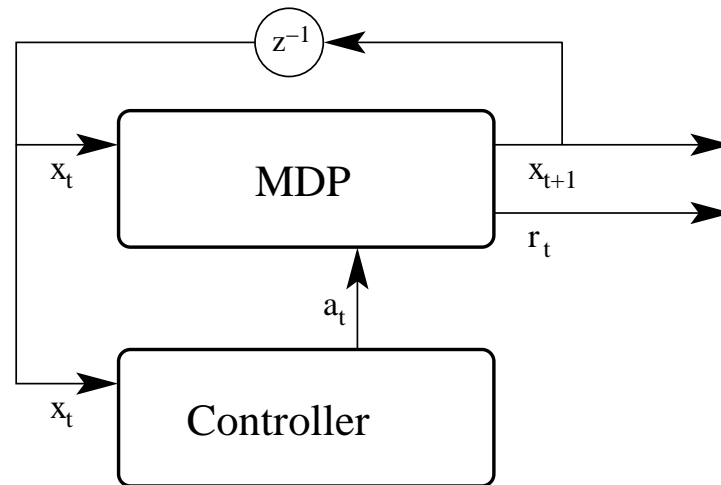
RL

EXAMPLE





MARKOV DECISION PROCESSES



State space:

\mathcal{X} , state $\mathbf{x} \in \mathcal{X}$

Action space:

\mathcal{A} , action $\mathbf{a} \in \mathcal{A}$

Joint state-action space:

$\mathcal{Z} = \mathcal{X} \times \mathcal{A}$, $\mathbf{z} = (\mathbf{x}, \mathbf{a})$

Transition prob. density:

$\mathbf{x}_{t+1} \sim p(\cdot | \mathbf{x}_t, \mathbf{a}_t)$

Reward prob. density:

$R(\mathbf{x}_t, \mathbf{a}_t) \sim q(\cdot | \mathbf{x}_t, \mathbf{a}_t)$



CONTROL AND RETURNS

Stationary policy:

$$\mathbf{a}_t \sim \mu(\cdot | \mathbf{x}_t)$$

Path:

$$\xi^\mu = (z_0, z_1, \dots)$$

Discounted Return:

$$D(\xi^\mu) = \sum_{i=0}^{\infty} \gamma^i R(z_i)$$

Value function:

$$V^\mu(\mathbf{x}) = \mathbf{E}_\mu[D(\xi^\mu) | \mathbf{x}_0 = \mathbf{x}]$$

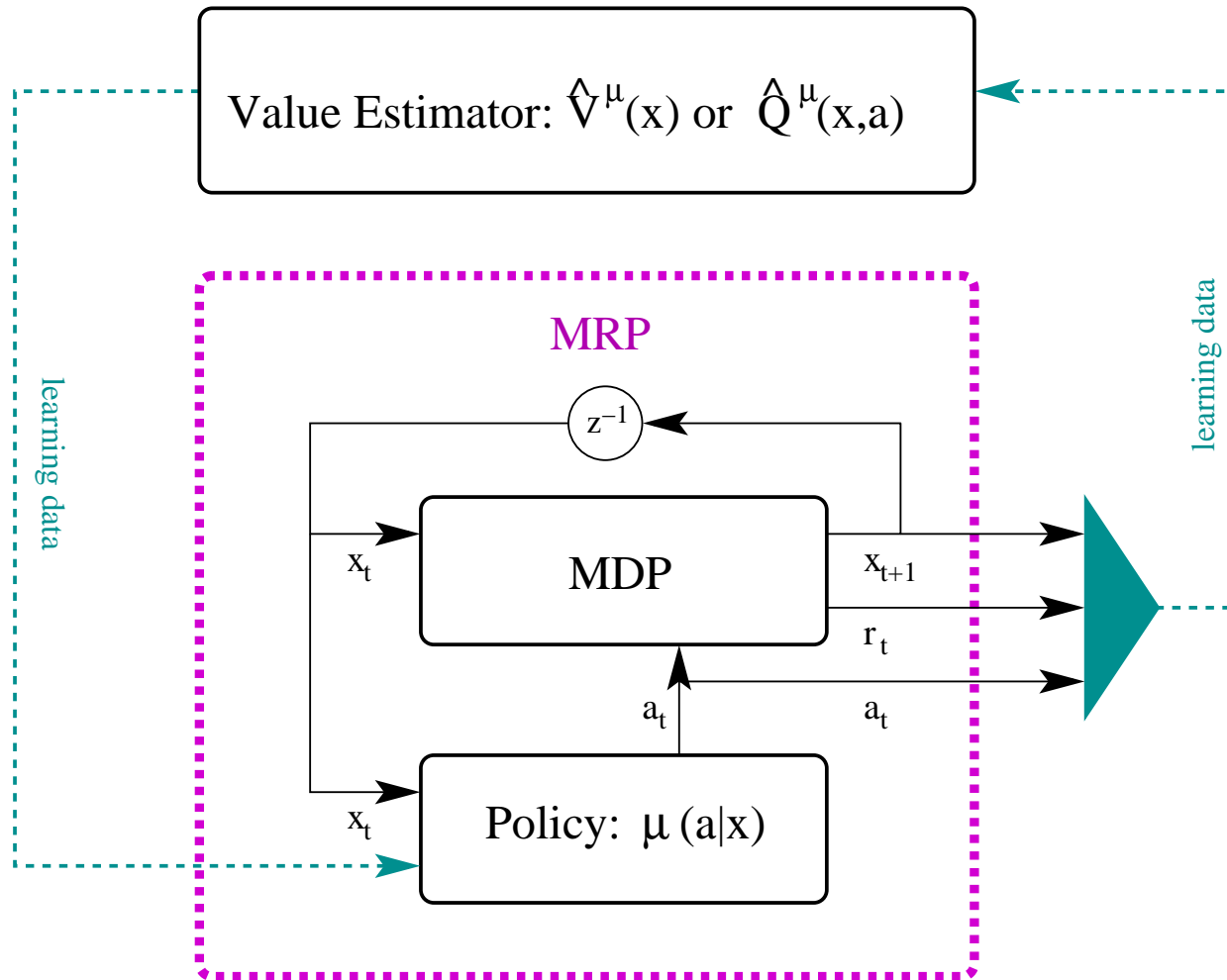
State-action value func.:

$$Q^\mu(z) = \mathbf{E}_\mu[D(\xi^\mu) | z_0 = z]$$

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

Note: If $Q^*(\mathbf{x}, \mathbf{a}) = Q^{\mu^*}(\mathbf{x}, \mathbf{a})$ is available, then an optimal action for state \mathbf{x} is given by any $\mathbf{a}^* \in \operatorname{argmax}_{\mathbf{a}} Q^*(\mathbf{x}, \mathbf{a})$.

VALUE-BASED RL



Bayes



RL

BELLMAN'S EQUATION

For a fixed policy μ :

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

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

RL Algorithms

Purely Policy based
(Policy Gradient)

Value-Function based

Value Iteration type
(Q-Learning)

Policy Iteration type
(Actor-Critic, OPI, SARSA)

PI methods need a “subroutine” for policy evaluation

Bayes



RL

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



GP 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, IID, 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:

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

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:

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

Bayes



RL

THE POSTERIOR

General noise covariance:

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

Joint distribution:

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

Condition on R_{t-1} :

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

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

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



LEARNING STATE-ACTION VALUES

Under a fixed stationary policy μ , state-action pairs \mathbf{z}_t form a Markov chain, just like the states \mathbf{x}_t .

Consequently $Q^\mu(\mathbf{z})$ behaves similarly to $V^\mu(\mathbf{x})$:

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

Posterior moments:

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

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



POLICY IMPROVEMENT

Optimistic Policy Iteration algorithms work by maintaining a policy evaluator \hat{Q}_t and selecting the action at time t semi-greedily w.r.t. to the current state-action value estimates $\hat{Q}_t(\mathbf{x}_t, \cdot)$.

Policy evaluator	Parameters	OPI algorithm
Online TD(λ) (Sutton)	\mathbf{w}_t	SARSA (Rummery & Niranjan)
Online GPTD (Engel et Al.)	α_t, \mathbf{C}_t	GPSARSA (Engel et Al.)



GPSARSA ALGORITHM

Initialize $\alpha_0 = \mathbf{0}$, $\mathbf{C}_0 = 0$, $\mathcal{D}_0 = \{z_0\}$, $\mathbf{c}_0 = \mathbf{0}$, $d_0 = 0$, $1/s_0 = 0$

for $t = 1, 2, \dots$

observe \mathbf{x}_{t-1} , \mathbf{a}_{t-1} , r_{t-1} , \mathbf{x}_t

$\mathbf{a}_t = \text{SemiGreedyAction}(\mathbf{x}_t, \mathcal{D}_{t-1}, \alpha_{t-1}, \mathbf{C}_{t-1})$

$d_t = \frac{\gamma \sigma_{t-1}^2}{s_{t-1}} d_{t-1} + \text{temporal difference}$

$\mathbf{c}_t = \dots$, $s_t = \dots$

$\alpha_t = \begin{pmatrix} \alpha_{t-1} \\ 0 \end{pmatrix} + \frac{\mathbf{c}_t}{s_t} d_t$

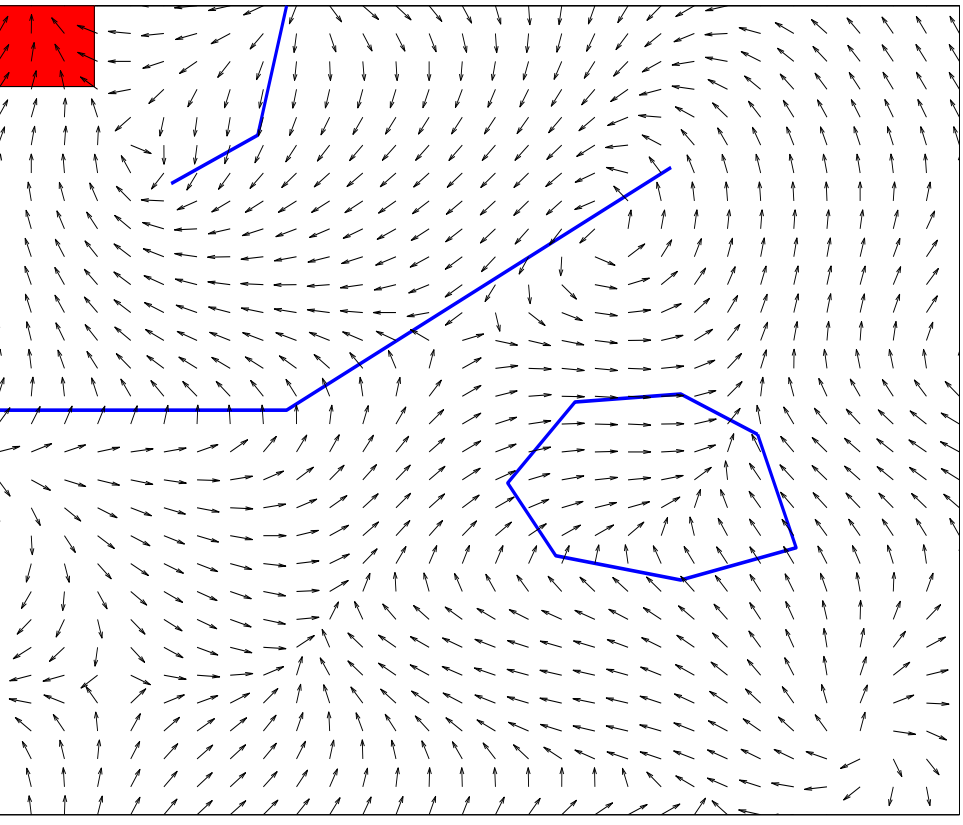
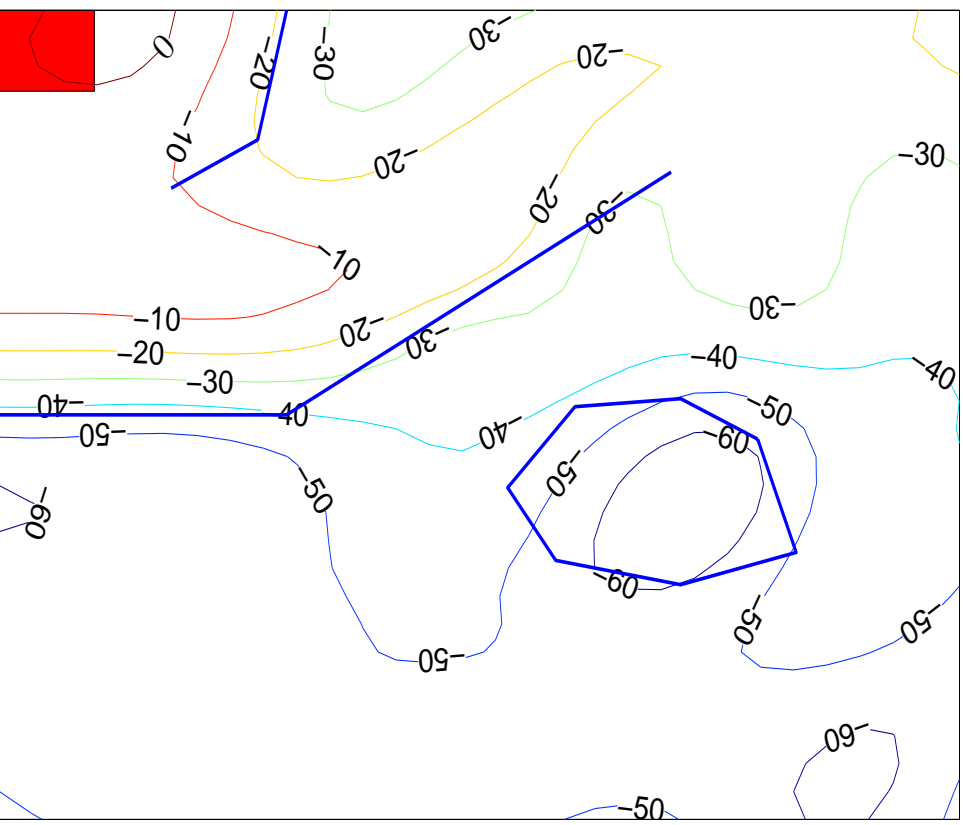
$\mathbf{C}_t = \begin{bmatrix} \mathbf{C}_{t-1} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix} + \frac{1}{s_t} \mathbf{c}_t \mathbf{c}_t^\top$

$\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{z_t\}$

end for

return α_t , \mathbf{C}_t , \mathcal{D}_t

A 2D NAVIGATION TASK



Bayes



RL

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?