# Bayesian Inference for Least-Squares Temporal Difference Regularization

## Nikolaos Tziortziotis[†] and Christos Dimitrakakis[⋆]

[†] DaSciM, LIX, École Polytechnique, France
[⋆] University of Lille, France - SEAS, Harvard University, USA
✉ email: {ntziorzi, christos.dimitrakakis}@gmail.com
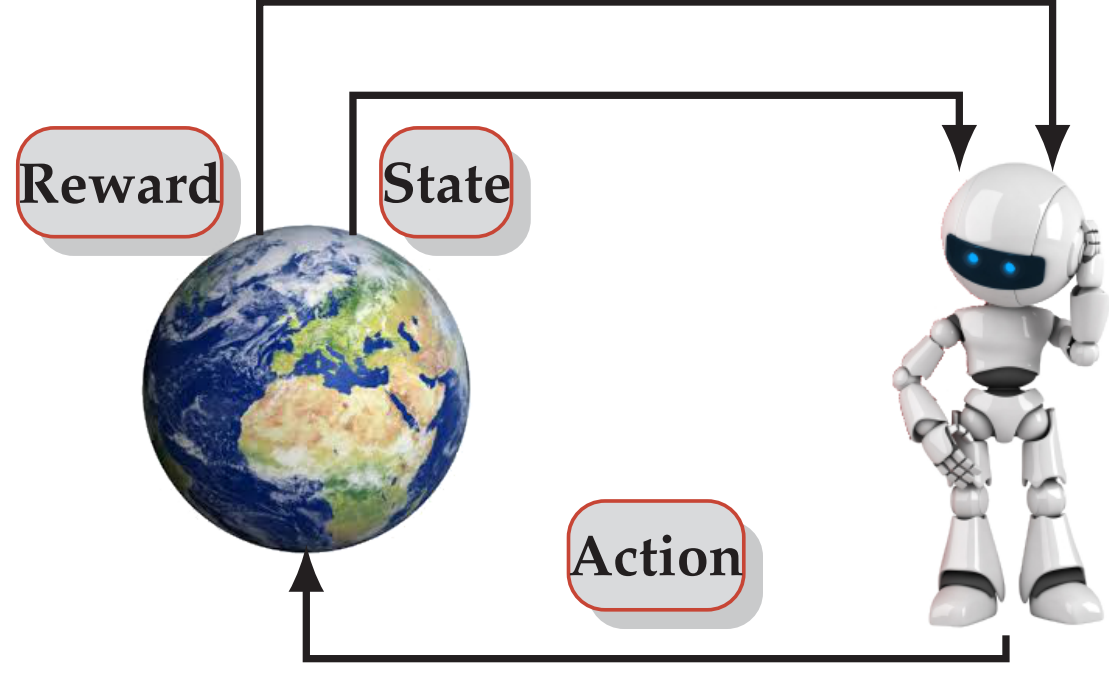
## Abstract

**Bayesian Least-Squares Temporal Difference:**
✓ Fully Bayesian approach for LSTD
✓ Probabilistic inference of value functions
✓ Quantifies our uncertainty about value function
**Variational Bayesian LSTD:**
✓ Sparse model - Good generalisation capabilities
✓ Automatically determine the model's complexity
✓ No need to select a regularization parameter

## Reinforcement Learning (RL)

Learning to act in an unknown environment, by
interaction and reinforcement.

RL tasks formulated as **MDPs**, $\{\mathcal{S}, \mathcal{A}, P, r, \gamma\}$
**Policy** $\pi$: $\mathcal{S} \rightarrow \mathcal{A}$ (map states to actions)
**Value function:** $V^\pi(s) \triangleq \mathbb{E}^\pi \left[ \sum_{t=0}^\infty \gamma^t r(s_t) | s_0 = s \right]$
**Bellman operator:**

$$(T^\pi V)(s) = r(s) + \gamma \int_\mathcal{S} V(s') dP(s'|s, \pi(s))$$

**V is the unique fixed-point of Bellman operator:**

$$V^\pi = T^\pi V^\pi \Rightarrow V^\pi = (I - \gamma P^\pi)^{-1} r$$

☺ The value function cannot be represented in an explicit way (continuous state space)
☺ The model of the MDP is unknown
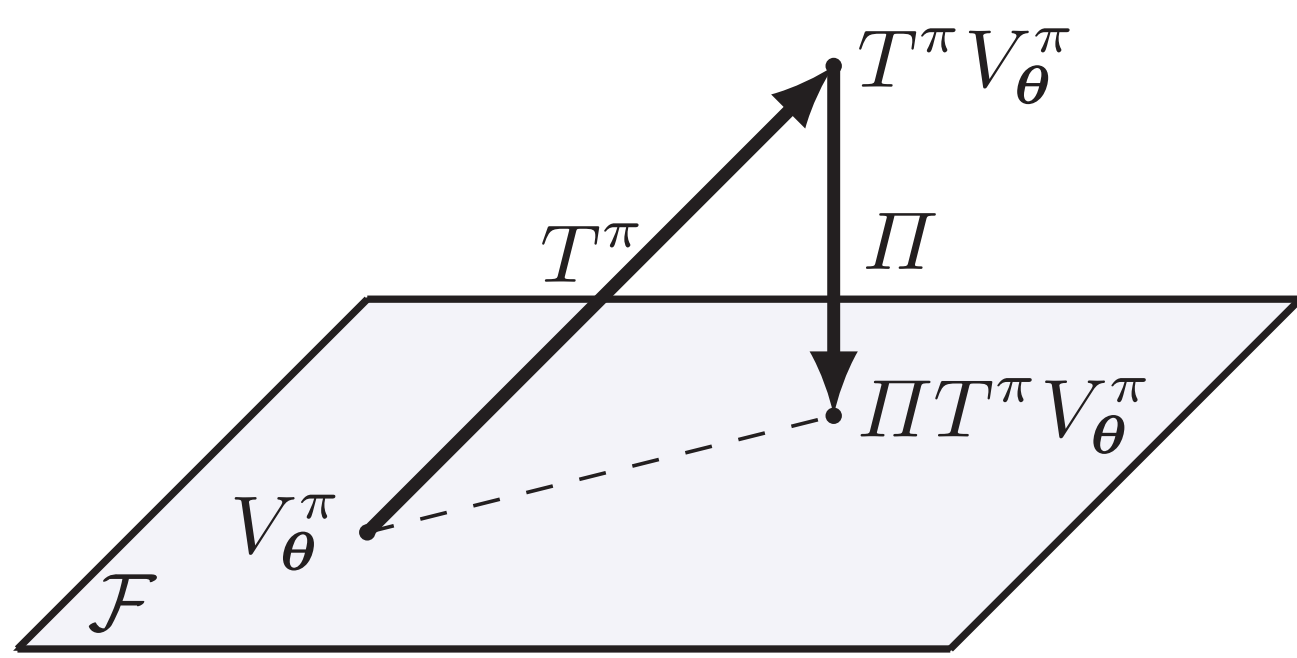
☺ **Value function approximation:**

$$V_\theta^\pi(s) = \phi(s)^\top \theta = \sum_{i=1}^k \phi_i(s) \theta_i,$$

where $\mathcal{F} = \{f_\theta | f_\theta(\cdot) = \phi(\cdot)^\top \theta\}$.
☺ **Access to a set of transitions:** $\mathcal{D} = \{(s_i, r_i, s_i')\}_{i=1}^n$,
where we define: $\tilde{\Phi} = [\phi(s_1)^\top; \dots; \phi(s_n)^\top]$,
$\tilde{R} = [r_1, \dots, r_n]^\top$ and $\tilde{\Phi}' = [\phi(s_1')^\top; \dots; \phi(s_n')^\top]$.

## Least-Squares Temporal Difference

Minimize the *mean-square projected Bellman error*
(MSPBE): $\theta = \arg\min_{\theta \in \mathbb{R}^k} \|V_\theta^\pi - \Pi T^\pi V_\theta^\pi\|_D^2$.

**Projection operator over** $\mathcal{F}$: $\Pi = \Phi C^{-1} \Phi^\top D$

### Nested Optimization Problem

$$\boldsymbol{u}^* = \arg\min_{\boldsymbol{u} \in \mathbb{R}^k} \|\Phi \boldsymbol{u} - T^\pi \Phi \theta\|_D^2 \quad \text{(Projection step)}$$

$$\theta = \arg\min_{\theta \in \mathbb{R}^k} \|\Phi \theta - \Phi \boldsymbol{u}^*\|_D^2 \quad \text{(Fixed-point step)}$$

### Solution

$$\boldsymbol{u}^* = \tilde{C}^{-1} \tilde{\Phi} (\tilde{R} + \gamma \tilde{\Phi}' \theta),$$

$$\theta = (\tilde{\Phi}^\top (\tilde{\Phi} - \gamma \tilde{\Phi}'))^{-1} \tilde{\Phi}^\top \tilde{R} = A^{-1} \boldsymbol{b},$$

where $\tilde{C} \triangleq \tilde{\Phi}^\top \tilde{\Phi}$, $A \triangleq \tilde{\Phi}^\top (\tilde{\Phi} - \gamma \tilde{\Phi}')$, and $\boldsymbol{b} \triangleq \tilde{\Phi}^\top \tilde{R}$.

As the number of samples $n$ increases, the LSTD solution $\tilde{\Phi}\theta$ converges to the fixed-point of $\hat{\Pi} T^\pi$.

### Regularized LSTD Schemes

- $\ell_2$-LSTD
- Lasso-TD
- LC-TD
- LARS-TD
- $\ell_1$-PBR
- $\ell_{2,2}$-LSTD
- $\ell_{2,1}$-LSTD
- Dantzig-LSTD
- ODDS-TD

## Bayesian LSTD

### Empirical Bellman operator:

$$\hat{T}^\pi V_\theta^\pi = \boldsymbol{r} + \gamma P^\pi V_\theta^\pi + N, \quad N \sim \mathcal{N}(0, \beta^{-1} I)$$

**Given observations** $\mathcal{D}$ **(LSTD solution):**

$$V_\theta^\pi = \hat{\Pi} \hat{T}^\pi V_\theta^\pi \Leftrightarrow \tilde{\Phi}^\top \tilde{R} = \tilde{\Phi}^\top (\tilde{\Phi} - \gamma \tilde{\Phi}') \theta + \tilde{\Phi}^\top N$$

**Linear regression model:** $\boldsymbol{b} = A\theta + \tilde{\Phi}^\top N$
**Likelihood function:** $p(\boldsymbol{b}|\theta, \beta) = \mathcal{N}(\boldsymbol{b}|A\theta, \beta^{-1}\tilde{C})$

✓ Maximum likelihood inference corresponds to standard LSTD solution

**Prior distribution over** $\theta$ : $p(\theta|\alpha) = \mathcal{N}(\theta|0, \alpha^{-1} I)$.
**Logarithm of posterior distribution:**

$$\ln p(\theta|\mathcal{D}) \propto -\frac{\beta}{2} E_\mathcal{D}(\theta) - \frac{\alpha}{2} \theta^\top \theta,$$

where $E_\mathcal{D}(\theta) = (\boldsymbol{b} - A\theta)^\top \tilde{C}^{-1} (\boldsymbol{b} - A\theta)$ (MSPBE).

✓ Maximum a posteriori inference corresponds to $\ell_2$ regularization ($\lambda = \alpha/\beta$)

**Posterior distribution:** $p(\theta|\mathcal{D}) = \mathcal{N}(\theta|\boldsymbol{m}, S)$

$$S = (\alpha I + \beta \underbrace{A^\top \tilde{C}^{-1} A}_{\Sigma})^{-1} \text{ and } \boldsymbol{m} = \beta S A^\top \tilde{C}^{-1} \boldsymbol{b}$$

**Predictive distribution:**

$$p(V_\theta^\pi(s^*)|s^*, \mathcal{D}) = \int_\theta p(V_\theta^\pi(s^*)|\theta, s^*) dp(\theta|\boldsymbol{b}, \alpha, \beta)$$
$$= \mathcal{N}(V_\theta^\pi(s^*)|\phi(s^*)^\top \boldsymbol{m}, \phi(s^*)^\top S \phi(s^*))$$

## Variational Bayesian LSTD

**Prior distribution:** $p(\theta|\alpha) = \prod_{i=1}^k \mathcal{N}(\theta_i|0, \alpha_i^{-1})$
**Hyperprior over** $\alpha$: $p(\alpha) = \prod_{i=1}^k Gamma(\alpha_i|h_a, h_b)$
**Hyperprior over** $\beta$: $p(\beta) = Gamma(\beta|h_c, h_d)$

**Posterior distribution over** $\mathcal{Z} = \{\theta, \alpha, \beta\}$:

$$p(\theta, \alpha, \beta|\boldsymbol{b}) = \frac{p(\boldsymbol{b}|\theta, \beta) p(\beta) p(\theta|\alpha) p(\alpha)}{p(\boldsymbol{b})}$$

☺ **Marginal likelihood is analytically intractable**

✓ **We resort to variational inference**

**Variational distribution:** $\mathcal{Q}(\mathcal{Z}) = \mathcal{Q}_\theta(\theta) \mathcal{Q}_\alpha(\alpha) \mathcal{Q}_\beta(\beta)$
**Optimal distribution for each factor:**

$$\mathcal{Q}_\theta(\theta) = \mathcal{N}(\theta|\boldsymbol{m}, S)$$
$$\mathcal{Q}_\beta(\beta) = Gamma(\beta|\tilde{c}, \tilde{d})$$
$$\mathcal{Q}_\alpha(\alpha) = \prod_{i=1}^k Gamma(\alpha_i|\tilde{a}_i, \tilde{b}_i)$$

where,

$$S = (diag \, \mathbb{E}[\alpha] + \mathbb{E}[\beta] \Sigma)^{-1}, \quad \boldsymbol{m} = \mathbb{E}[\beta] S A^\top \tilde{C}^{-1} \boldsymbol{b},$$
$$\tilde{a}_i = h_a + \frac{1}{2}, \quad \tilde{b}_i = h_b + \frac{1}{2} \mathbb{E}[\theta_i^2],$$
$$\tilde{c} = h_c + \frac{k}{2}, \quad \tilde{d} = h_d + \frac{1}{2} \|\boldsymbol{b} - A\boldsymbol{m}\|_{\tilde{C}}^2 + \frac{1}{2} tr(\Sigma S).$$

### Lower bound

$$\mathcal{L}(\mathcal{Q}) = \frac{1}{2} \ln |S| - \frac{1}{2} |\tilde{C}| + \sum_{i=1}^k \{\ln \Gamma(\tilde{a}_i) - \tilde{a}_i \ln \tilde{b}_i\}$$
$$+ \ln \Gamma(\tilde{c}) - \tilde{c} \ln \tilde{d} + \frac{k}{2}(1 - \ln 2\pi) - k \ln \Gamma(h_a)$$
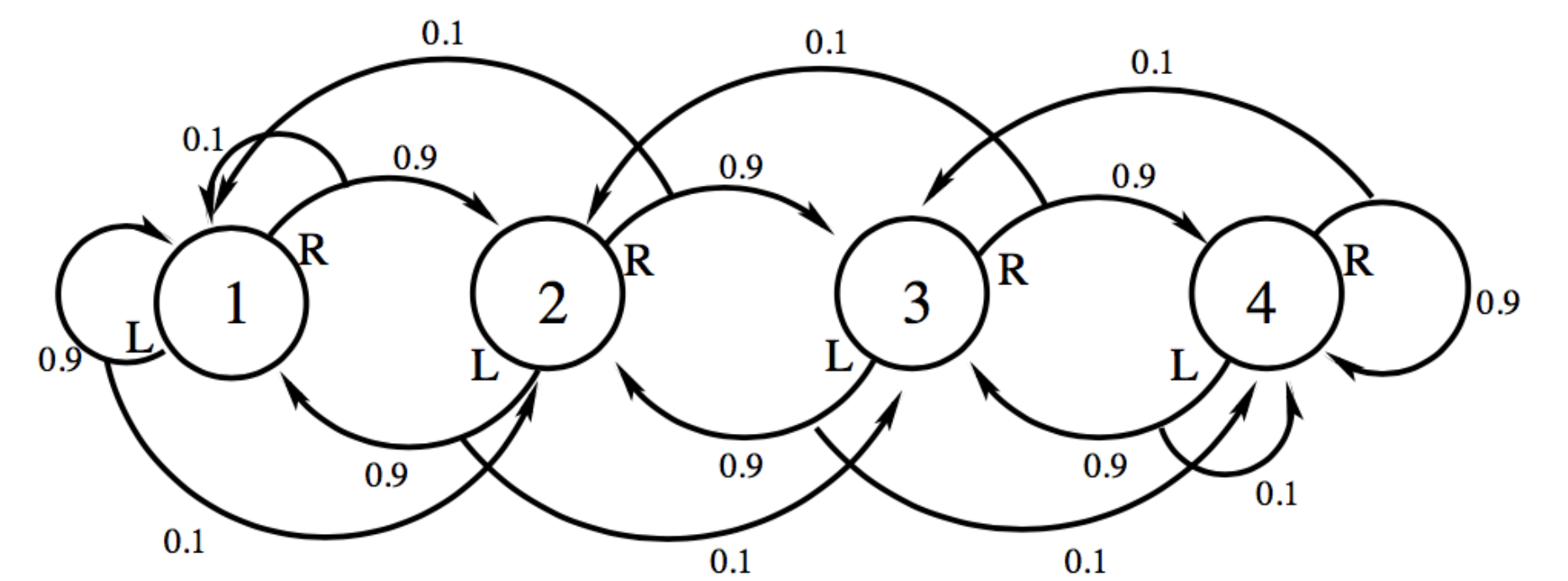$$+ k h_a \ln h_b - \ln \Gamma(h_c) + h_c \ln h_d$$

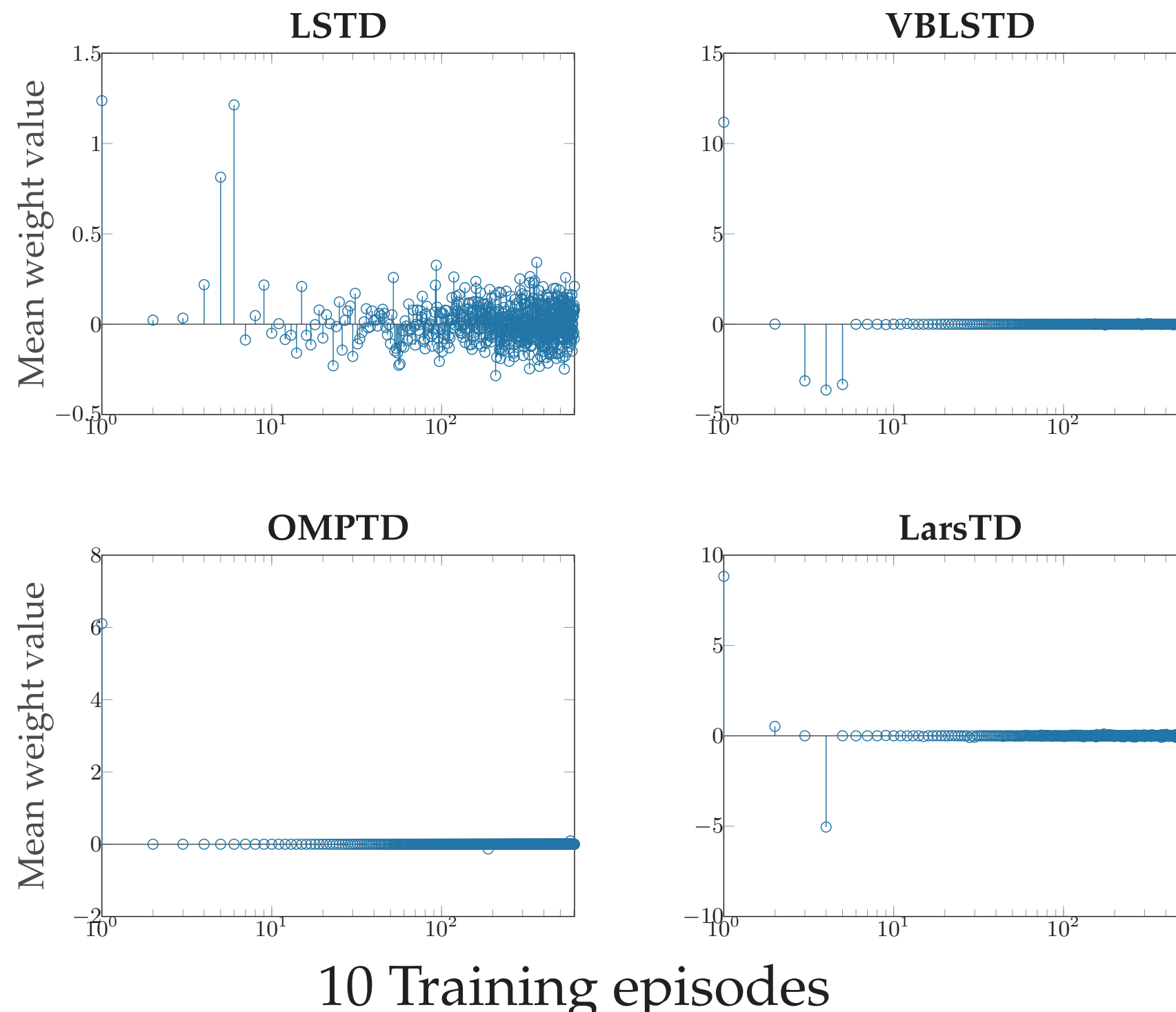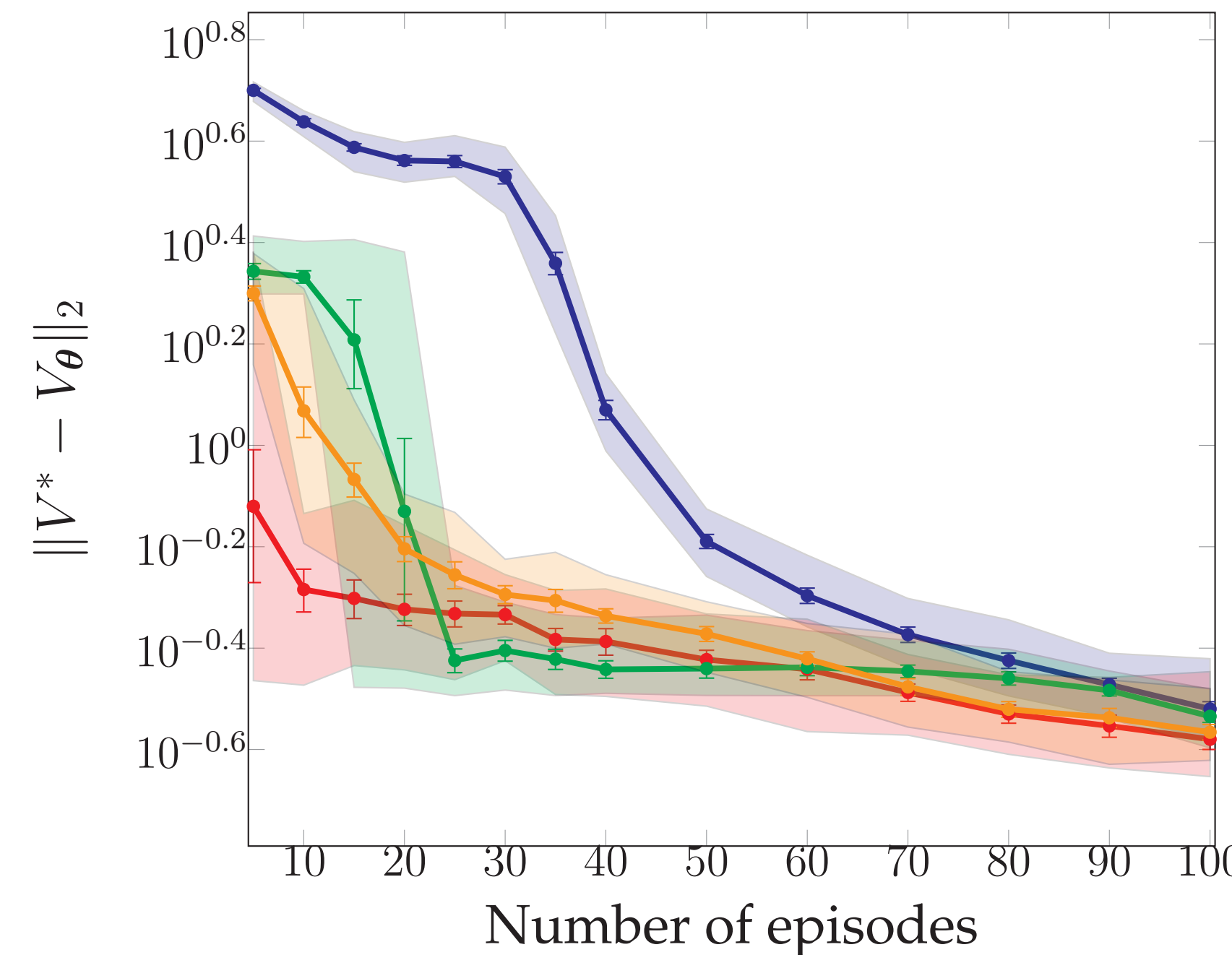## Experimental Results

### Corrupted Chain

- A 20-state, 2-actions (*left* or *right*) MDP
- The probability of action's success is equal to 0.9
- A reward of 1 is given only at the ends of chain
- The horizon of each episode is set equal to 20
- Optimal policy: i) $s \leq 10$: a = L, ii) $s > 10$: a = R
- Value function representation:

$$\phi(s) = (1, \text{RBF}_1(s), \dots, \text{RBF}_5(s), x_1, \dots, x_{600})$$

where $x_i \sim \mathcal{N}(0, 1), \forall i \in \{1, \dots, 600\}$
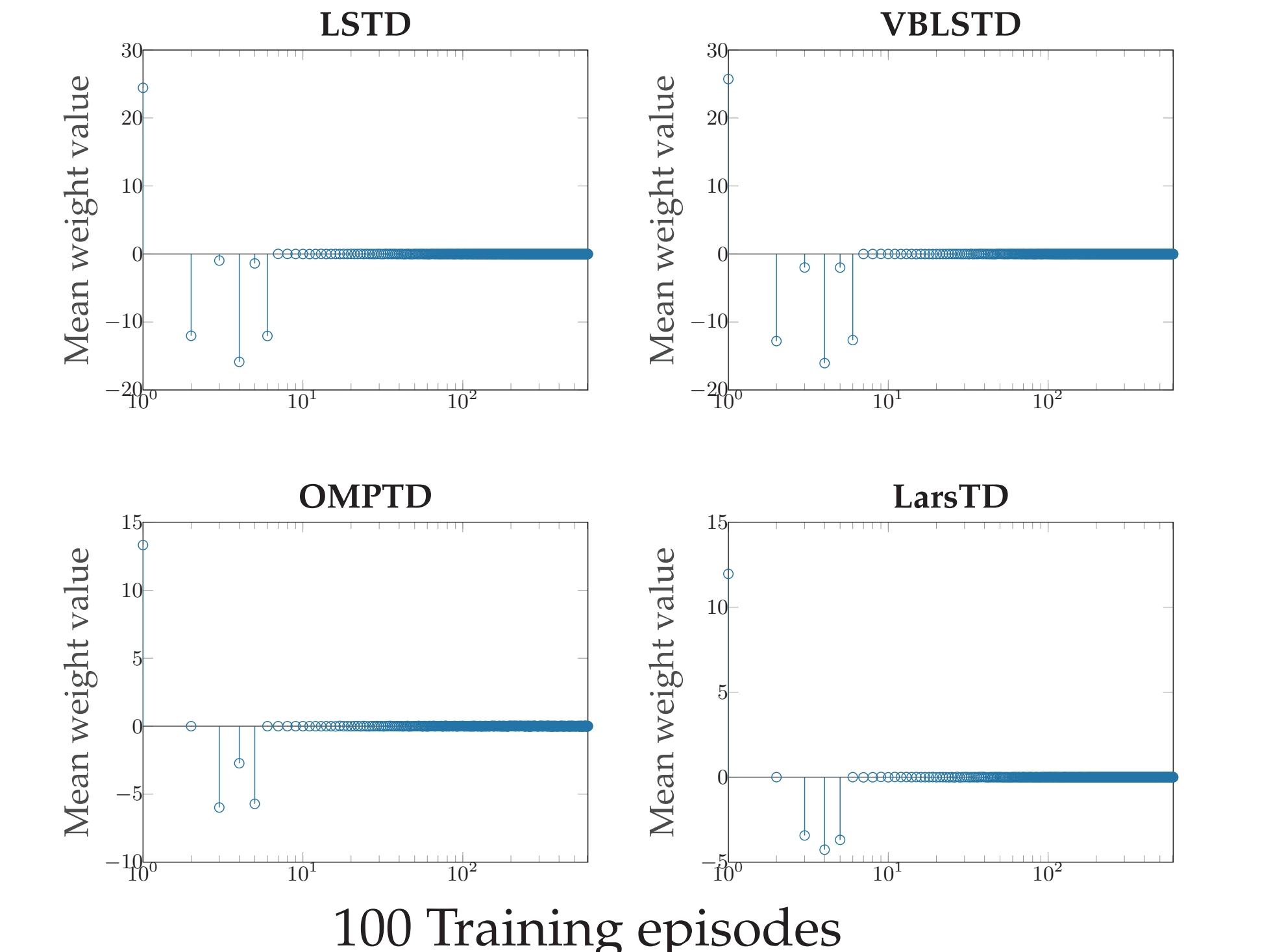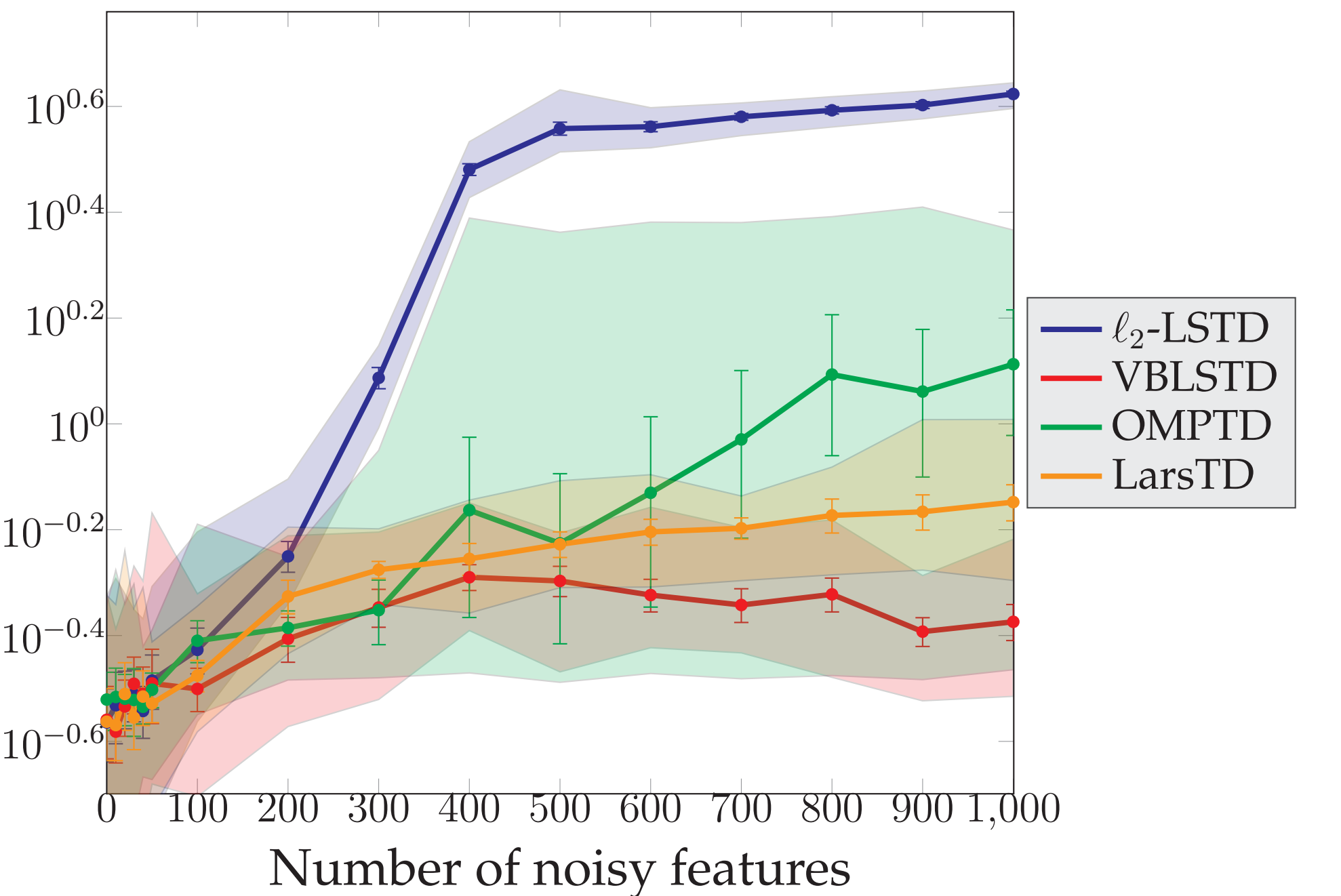




Corrupted chain varying the number of samples

Corrupted chain varying the number of noise features

$\ell_2$-LSTD
VBLSTD
OMPTD
LarsTD



10 Training episodes

100 Training episodes

The 606 mean weight values. The first weight is the bias term, the next 5 correspond to the relevant features (RBFs), and the rest 600 correspond to the noise (irrelevant) features.