# Incremental Least-Square Temporal Difference Learning (iLSTD)

## Alborz Geramifard, Michael Bowling, Richard S. Sutton
## University of Alberta

## Problem

- Markov Decision Process (MDP)

$$(\mathcal{S}, \mathcal{A}, \mathcal{P}_{ss'}^a, \mathcal{R}_{ss'}^a, \gamma)$$

- We focus on online policy evaluation

$$V^\pi(s) = \sum_a \pi(s,a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

## Related Work

- Temporal Difference Learning: TD(0)

TD

S $\xrightarrow{(\pi)\ a,r}$ S'

$$\delta_t(V) = r_{t+1} + \gamma V(s_{t+1}) - V(s_t).$$

- Using Linear Function Approximation

$$V(s_t) = \phi(s_t)^T \boldsymbol{\theta}$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha_t \phi(s_t) \delta_t(V)$$

we assume $k$ features are "on" on each time step

- Least-Square TD (LSTD)

LSTD

$$\boldsymbol{\mu}_t(\boldsymbol{\theta}) = \sum_{i=1}^t \phi_i \delta_i(V_{\boldsymbol{\theta}})$$

$$= \underbrace{\sum_{i=1}^t \phi_i r_{i+1}}_{\mathbf{b}_t} - \underbrace{\sum_{i=1}^t \phi_i (\phi_i - \gamma\phi_{i+1})^T \boldsymbol{\theta}}_{\mathbf{A}_t}$$

$$\boldsymbol{\theta}_{t+1} = \mathbf{A}_t^{-1} \mathbf{b}_t.$$

## New Approach

- TD
  - Cheap - O(k) per time-step
  - Relatively Data Inefficient

- LSTD
  - Expensive - O(n^2) per time-step
  - Data Efficient

- iLSTD
  - Cheap - O(mn+k^2) per time-step
  - Data Efficient

iLSTD

```
0   s ← s_0, A ← 0, μ ← 0, t ← 0
1   Initialize θ arbitrarily
2   repeat
3       Take action according to π and observe r, s'
4       t ← t + 1
5       Δb ← φ(s)r
6       ΔA ← φ(s)(φ(s) − γφ(s'))^T
7       A ← A + ΔA
8       μ ← μ + Δb − (ΔA)θ
9       for i from 1 to m do
10          j ← argmax(|μ_j|)
11          θ_j ← θ_j + αμ_j
12          μ ← μ − αμ_j A e_i
13      end for
14  end repeat
```

- Would like to Update **θ** by the sum of the TD updates (**μ**)
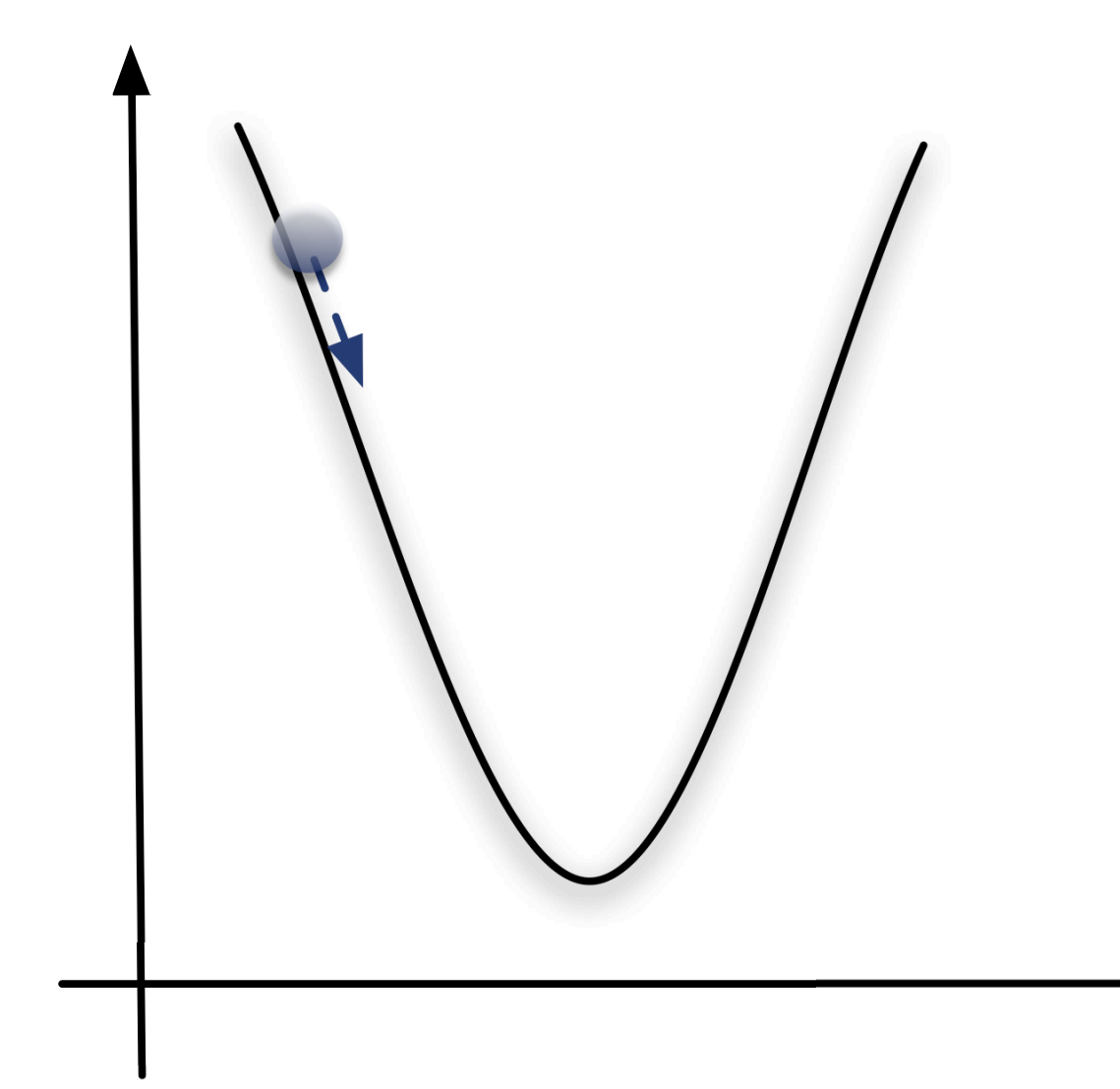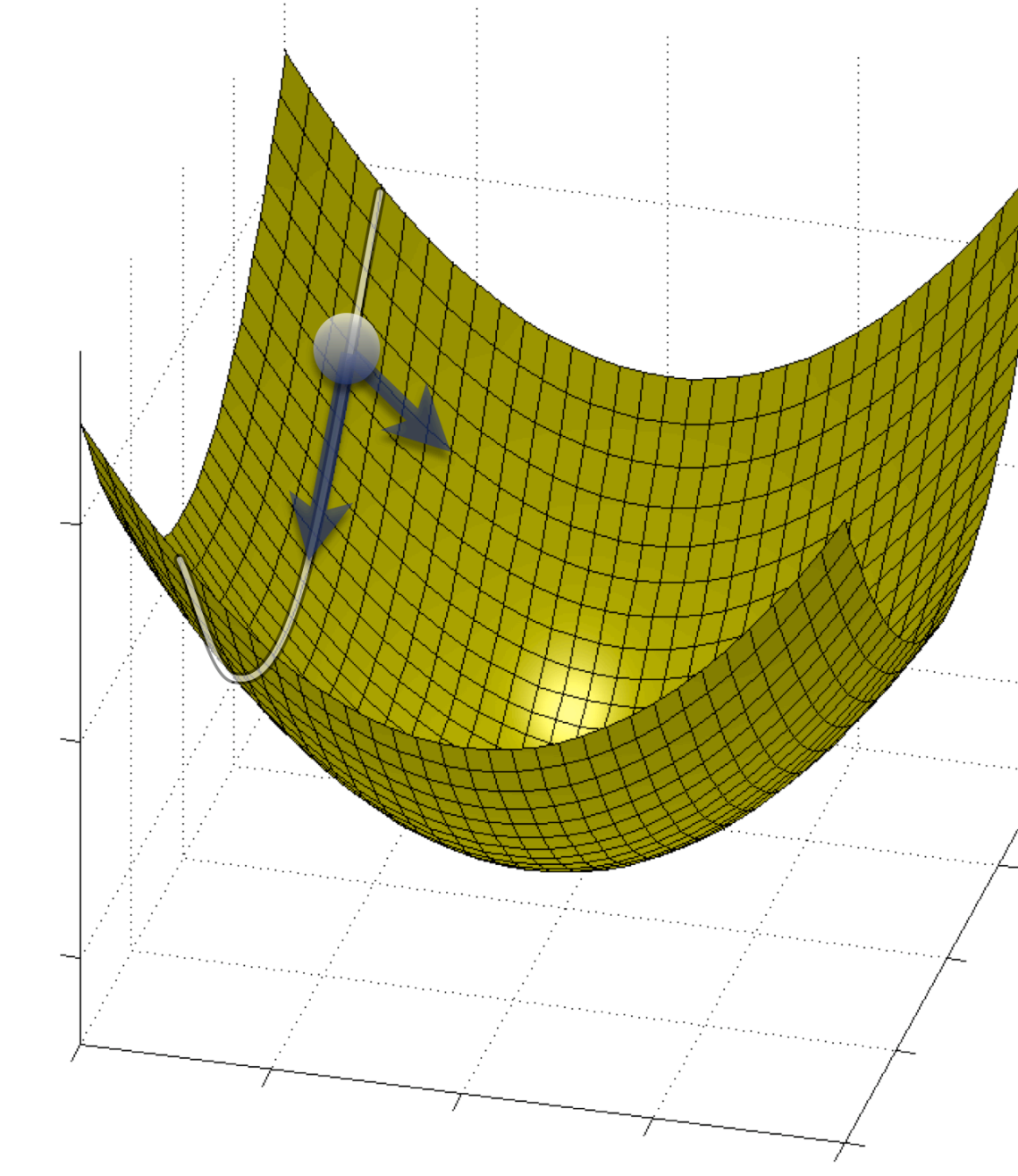- Pick the dimension with the largest TD update
- Descend in that dimension

- Updates:

$$\boldsymbol{\mu}_t(\boldsymbol{\theta}_{t+1}) = \boldsymbol{\mu}_t(\boldsymbol{\theta}_t) - \mathbf{A}_t(\Delta\boldsymbol{\theta}_t)$$

$$\boldsymbol{\mu}_t(\boldsymbol{\theta}_t) = \boldsymbol{\mu}_{t-1}(\boldsymbol{\theta}_t) + \Delta\mathbf{b}_t - (\Delta\mathbf{A}_t)\boldsymbol{\theta}_t$$

$$O(mn + k^2)$$

- Maximum number of "on" features
- Number of features
- Number of descents per time-step

## Results

Boyan's MDP

26 Features



101 Features



Timing

AAAI - 2006 Boston