

Temporal Difference Bayesian Model Averaging

A Bayesian Perspective on Adapting λ

Carlton Downey (downeycarl@ecs.vuw.ac.nz)

Scott Sanner (ssanner@nicta.com.au)

TD(λ) - Motivation

- TD(λ) is one of the most widely used reinforcement learning algorithms in the world. However...
 - **Question:** data-dependent λ adaptation?
 - no manual λ tuning
 - adapt λ online
 - **Answer:** Bayesian model averaging (BMA)
 - what model?
 - how to compute efficiently online?
-

Outline

- Motivation
 - TD(λ) and Bayesian Model Averaging
 - Derivation
 - Comparative evaluation
 - Conclusions & future work
-

Markov Decision Process (MDP)

- MDP is a tuple $\langle \mathbf{S}, \mathbf{A}, \mathbf{T}, \mathbf{R} \rangle$
 - \mathbf{S} : finite set of states
 - \mathbf{A} : finite set of actions
 - $\mathbf{T}(\mathbf{s}', \mathbf{a}, \mathbf{s}) = \mathbf{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a})$: transition function
 - stationary, Markovian
 - $\mathbf{R}(\mathbf{s}, \mathbf{a})$: (stochastic) reward function
-

MDP Optimization

- γ ($0 \leq \gamma < 1$): discount factor
- $\pi(\mathbf{s}, \mathbf{a})$: (stochastic) exploration policy
 - $\pi(s, a) = P(a|s)$
- Objective to optimize:

$$Q_{\pi}(s, a) = E_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \cdot r_{t+1} \mid s_0 = s, a_0 = a \right]$$

TD

- TD Update:

$$Q_{\pi}(s_t, a_t) = Q_{\pi}(s_t, a_t) + \alpha [\Delta Q_t^{\lambda}]$$

$$\Delta Q_t^{\lambda} = R_t - Q_{\pi}(s_t, a_t)$$

- An n-step return is a bootstrapped estimate.

$$R_t^{(n)} = \sum_{i=1}^n \gamma^{i-1} r_{t+i} + \gamma^n Q_{\pi}(s_{t+n}, a_{t+n})$$

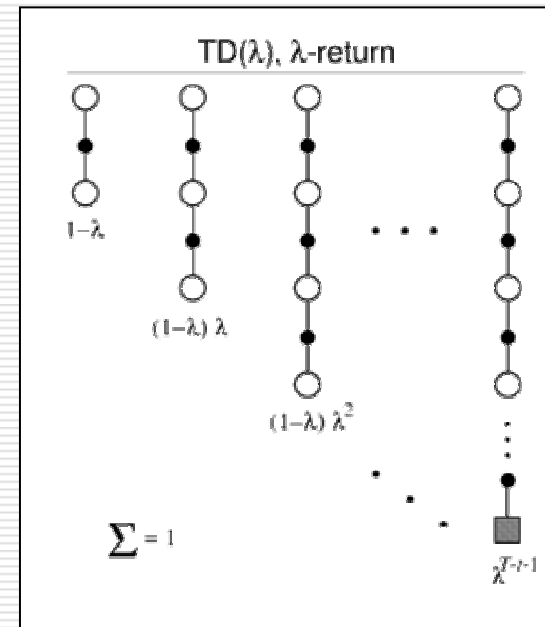
TD(λ)

- Averaging many n-step returns gives us TD(λ).

$$R_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} R_t^{(n)}$$

- Expressed recursively

$$R_t^\lambda = r_{t+1} + \gamma \left[(1 - \lambda) Q_\pi(s_{t+1}, a_{t+1}) + \lambda R_{t+1}^\lambda \right]$$



SARSA(λ)

- We thus obtain the TD(λ) update rule SARSA(λ):

$$Q_{\pi}(s_t, a_t) = Q_{\pi}(s_t, a_t) + \alpha[\Delta Q_t^{\lambda}]$$

$$\Delta Q_t^{\lambda} = R_t^{\lambda} - Q_{\pi}(s_t, a_t)$$

What's wrong with TD(λ)?

- ❑ TD(λ) chooses a particular fixed function of λ .
 - ❑ Is there a better way to weight each n-step return?
 - ❑ Can we weight them with so as to reduce variance?
-

BMA view of TD(λ)

- ❑ BMA used for reducing estimator variance.
 - ❑ BMA provides an intuitive data-dependent way to adjust the weights of multiple estimators.
 - ❑ BMA weights models according to how much the data supports them.
-

BMA view of TD(λ)

- Can we look at TD(λ) from BMA perspective?
- We derive an expected return model R_t^{BMA} :

$$\begin{aligned} R_t^{\text{BMA}} &= E_{P(\vec{q}|D)}[q_{s,a}] \\ &= \int_{q_{s,a} \in R} q_{s,a} P(q_{s,a} | D) dq_{s,a} \\ &= \int_{q_{s,a} \in R} q_{s,a} \sum_{m \in MC, TD} P(q_{s,a} | m) p(m | D) dq_{s,a} \\ &= \boxed{p(TD | D) E_{P(q_{s,a}|TD)}[q_{s,a}] + p(MC | D) E_{P(q_{s,a}|MC)}[q_{s,a}]} \end{aligned}$$

BMA view of TD(λ)

- If we set $P(\text{MC}|D) = \lambda$ and $P(\text{TD}|D) = (1 - \lambda)$:

$$\begin{aligned} R_t^{BMA} &= E_{P(q_{s,a}|\text{TD})}[q_{s,a}]p(\text{TD}|D) + E_{P(q_{s,a}|\text{MC})}[q_{s,a}]p(\text{MC}|D) \\ &= E_{P(q_{s,a}|\text{TD})}[q_{s,a}](1 - \lambda) + E_{P(q_{s,a}|\text{MC})}[q_{s,a}]\lambda \\ &= R_t^\lambda \end{aligned}$$

- We have exactly re-derived SARSA(λ).
 - But from a BMA perspective should we fix λ ?
 - We have data, we want to make lambda data dependent.
-

Gaussian Case

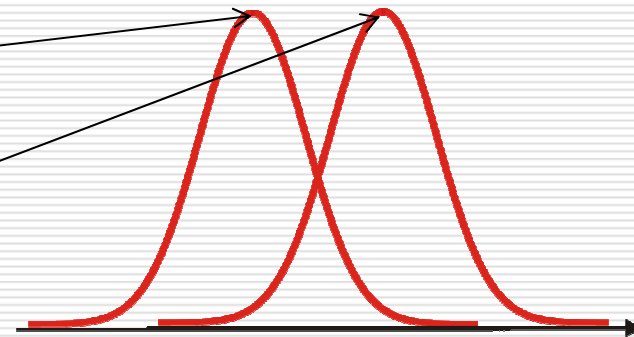
- Assume value of each pair (s,a) is independently gaussian.

- $P(q_{s,a} | m) = \mathcal{N}(q_{s,a}; \mu_{s,a}^m, (\sigma_{s,a}^m)^2)$

- For each model Mean and S.D. are sufficient statistics.

- $\mu_{s,a}^{MC} = q_{s,a}^{MC} = r_{t+1} + \gamma \mathcal{R}_{t+1}^\lambda$

$$\mu_{s,a}^{TD} = q_{s,a}^{TD} = r_{t+1} + \gamma Q_\pi(s_{t+1}, a_{t+1})$$



- To avoid costly computation we use the S.D of D.
-

Gaussian Case

□ Model Prediction is trivial:

□ Model Weight is harder:

$$E_{P(q_{s,a}|m)} [q_{s,a}] = \mu_{s,a}^m = q_{s,a}^m$$

$$\begin{aligned}
 P(m|D) &= P(m|D_{s,a}) \propto P(D_{s,a}|m)P(m) && \text{Bayes} \\
 &\propto P(D_{s,a}|m) && \text{Indep.} \\
 &= \prod_{d \in D_{s,a}} P(d|m) && \text{Unif. Prior} \\
 &= \prod_{d \in D_{s,a}} N(d; q_{s,a}^m, \sigma_{s,a}^2) && \text{i.i.d} \\
 &= \prod_{d \in D_{s,a}} C e^{-\frac{(d - q_{s,a}^m)^2}{2\sigma_{s,a}^2}} \\
 &= C^{|D_{s,a}|} e^{\left(\frac{1}{2\sigma_{s,a}^2}\right) \sum_{d \in D_{s,a}} (-d^2 - (q_{s,a}^m)^2 + 2q_{s,a}^m d)}
 \end{aligned}$$

Gaussian Case

- Simplifying exponent of model weight further...

$$\sum_{d \in D_{s,a}} (-d^2 - (q_{s,a}^m)^2 + 2q_{s,a}^m d)$$

$$= -\sum_{d \in D_{s,a}} d^2 - \sum_{d \in D_{s,a}} (q_{s,a}^m)^2 + 2q_{s,a}^m \sum_{d \in D_{s,a}} d$$

$$= -|D_{s,a}|(\sigma_{s,a}^2 + (\mu_{s,a}^D)^2) - |D_{s,a}|(q_{s,a}^m)^2 + 2q_{s,a}^m |D_{s,a}| \mu_{s,a}^D$$

Simplified
using def. of
variance

$$= -|D_{s,a}|(\sigma_{s,a}^2 + (\mu_{s,a}^D)^2 + (q_{s,a}^m)^2 - 2q_{s,a}^m \mu_{s,a}^D)$$

$$= -|D_{s,a}| \sigma_{s,a}^2 - |D_{s,a}| (q_{s,a}^m - \mu_{s,a}^D)^2$$

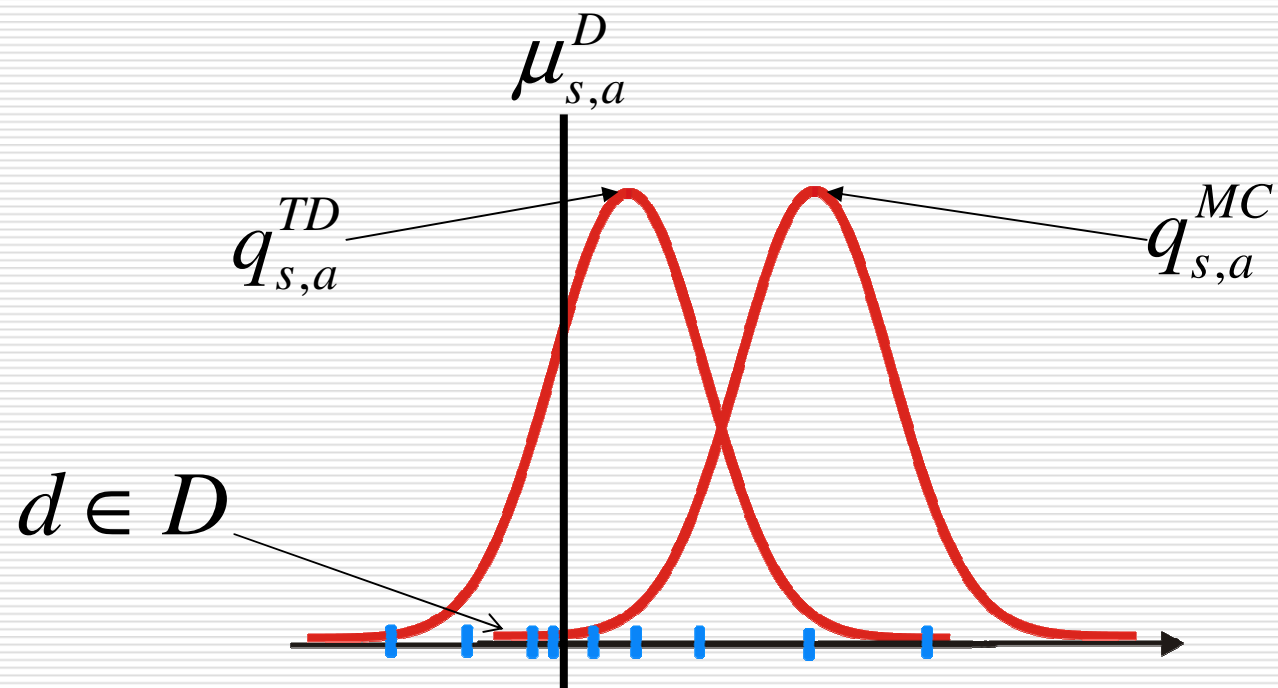
Gaussian Case

- Finally, substituting back for ...

$$\begin{aligned} P(m | D) &\propto C^{|D_{s,a}|} e^{-\frac{|D_{s,a}| \sigma_{s,a}^2}{2\sigma_{s,a}^2} - \frac{|D_{s,a}|}{2\sigma_{s,a}^2} (q_{s,a}^m - \mu_{s,a}^D)^2} \\ &= e^{-\frac{|D_{s,a}|}{2}} \left[N(q_{s,a}^m; \mu_{s,a}^D, \sigma_{s,a}^2) \right]^{|D_{s,a}|} \end{aligned}$$

- This result depends only $|D_{s,a}|$, $\mu_{s,a}^m$ and $\sigma_{s,a}^m$.
 - These can be calculated online in constant time.
-

Intuition



□ Clearly in this case $P(TD_{s,a} | D) > P(MC_{s,a} | D)$

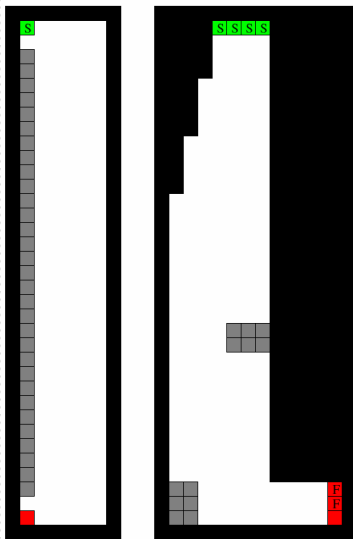
Algorithm

$$R_t^\lambda = r_{t+1} + \gamma[(1 - \lambda)Q_\pi(s_{t+1}, a_{t+1}) + \lambda R_{t+1}^\lambda]$$

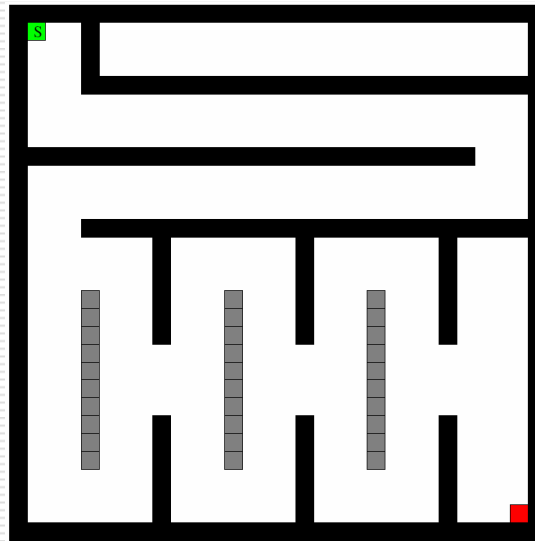
$$\begin{aligned}\lambda_{BMA} &= \frac{P(MC_{s,a} \mid D)}{P(MC_{s,a} \mid D) + P(TD_{s,a} \mid D)} \\ &= \frac{\left[N(q_{s,a}^{MC}; \mu_{s,a}^D, (\sigma_{s,a}^D)^2) \right]^{D_{s,a}|}}{\left[N(q_{s,a}^{MC}; \mu_{s,a}^D, (\sigma_{s,a}^D)^2) \right]^{D_{s,a}|} \left[N(q_{s,a}^{TD}; \mu_{s,a}^D, (\sigma_{s,a}^D)^2) \right]^{D_{s,a}|}}\end{aligned}$$

$$R_t^{BMA} = r_{t+1} + \gamma[(1 - \lambda_{BMA})Q_\pi(s_{t+1}, a_{t+1}) + \lambda_{BMA} R_{t+1}^{BMA}]$$

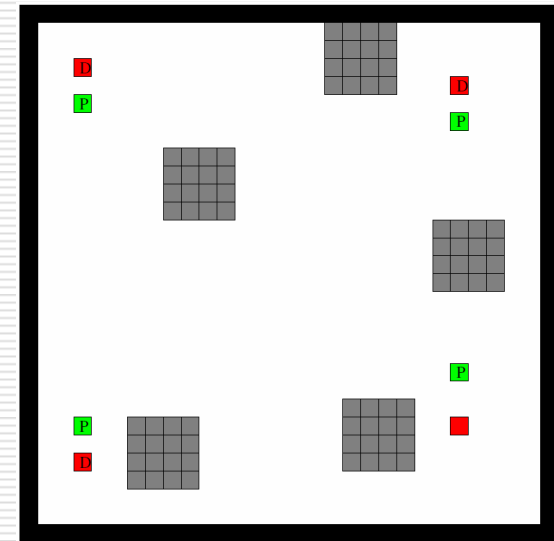
Experimental Setup



(a) Cliff (b) Corner

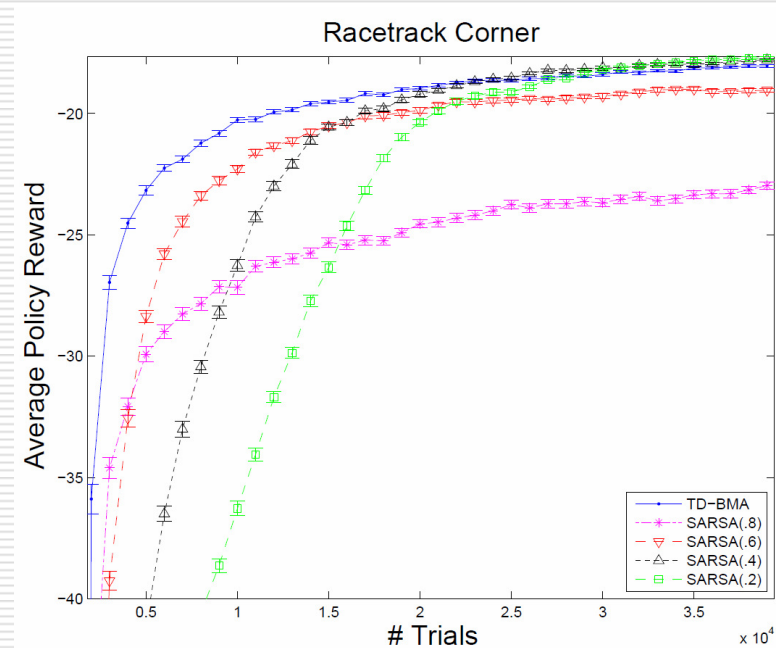
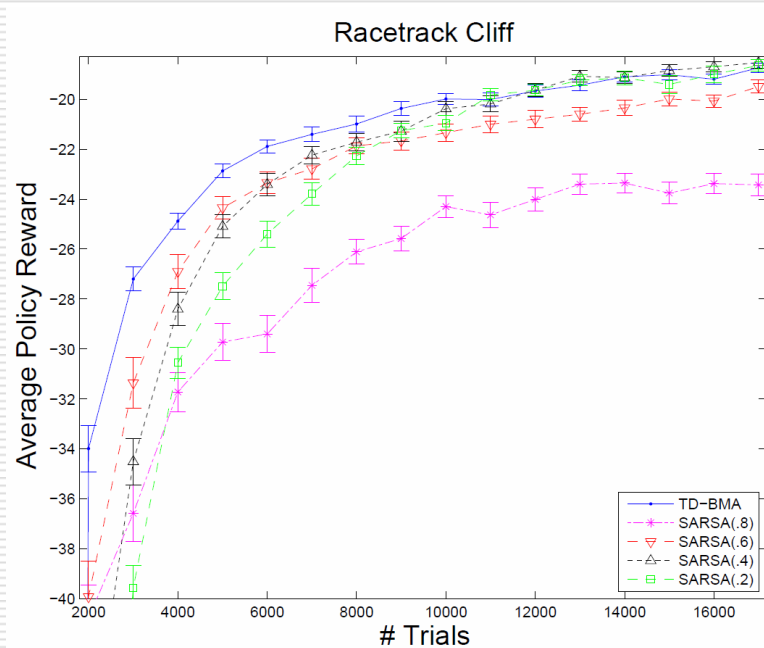


(c) Curves



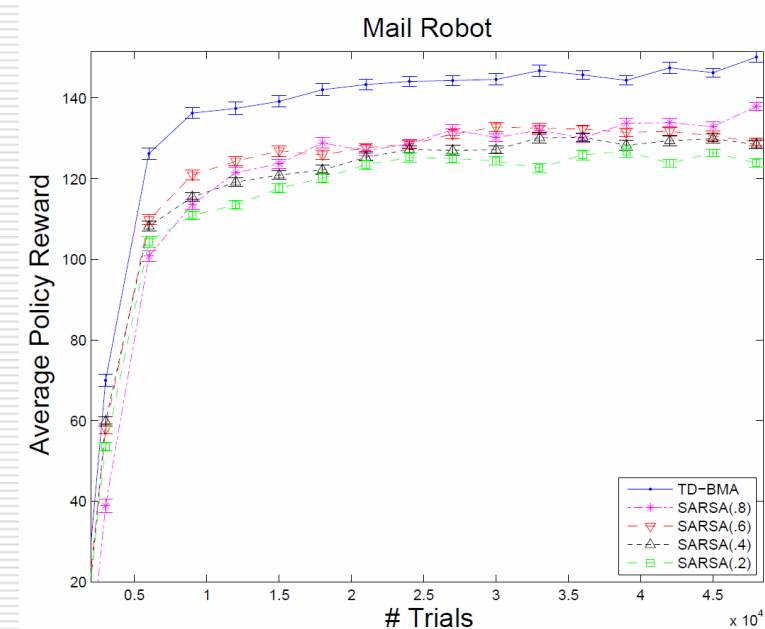
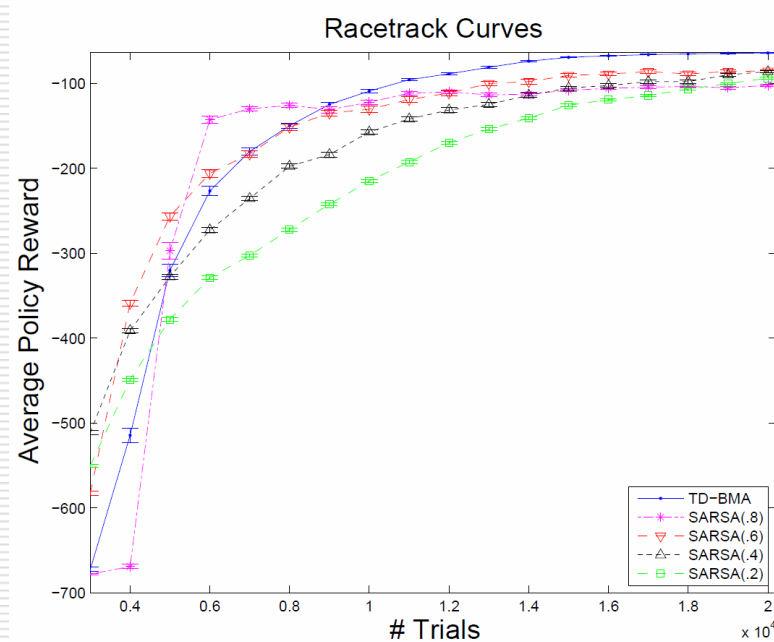
(d) Mail

Results: Varying λ



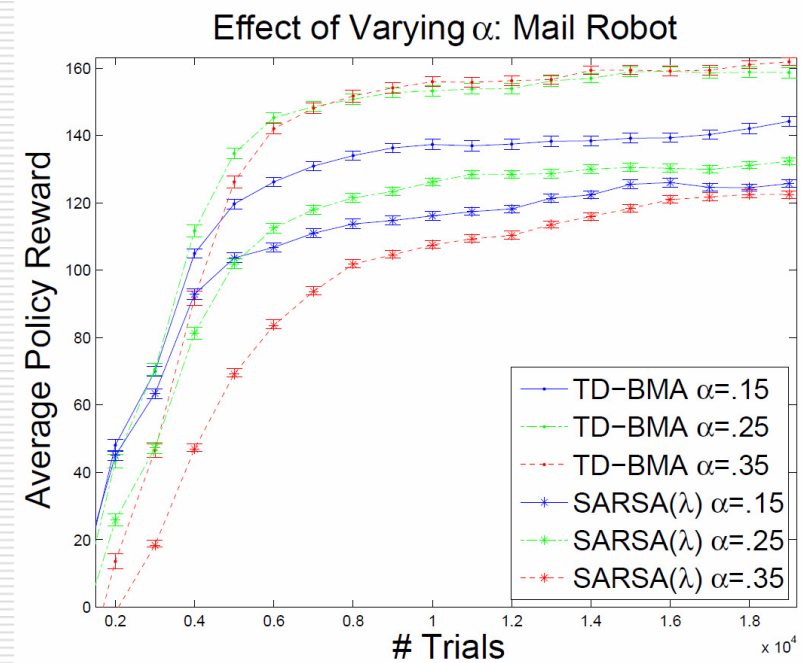
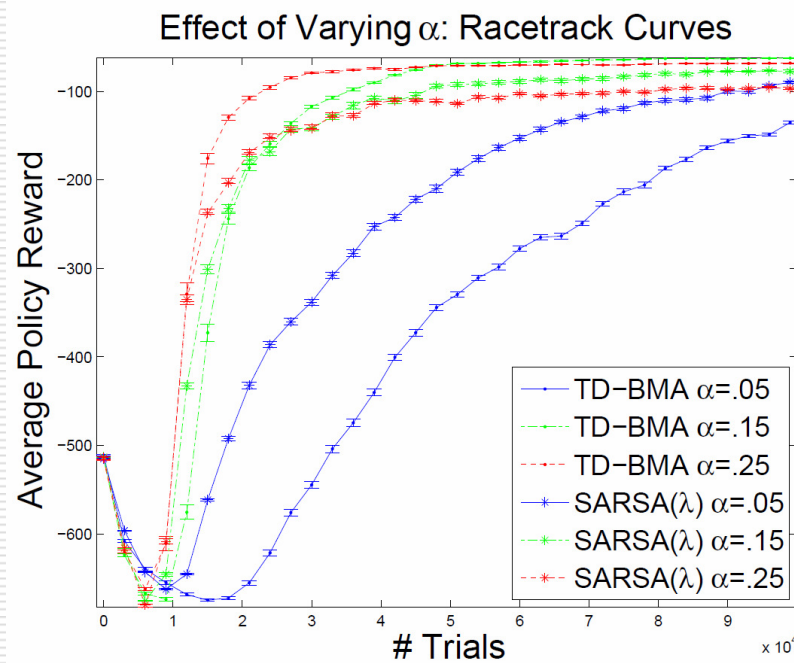
□ TD-BMA outperforms SARSA for all values of λ .

Results: Varying λ



□ TD-BMA outperforms SARSA for all values of λ .

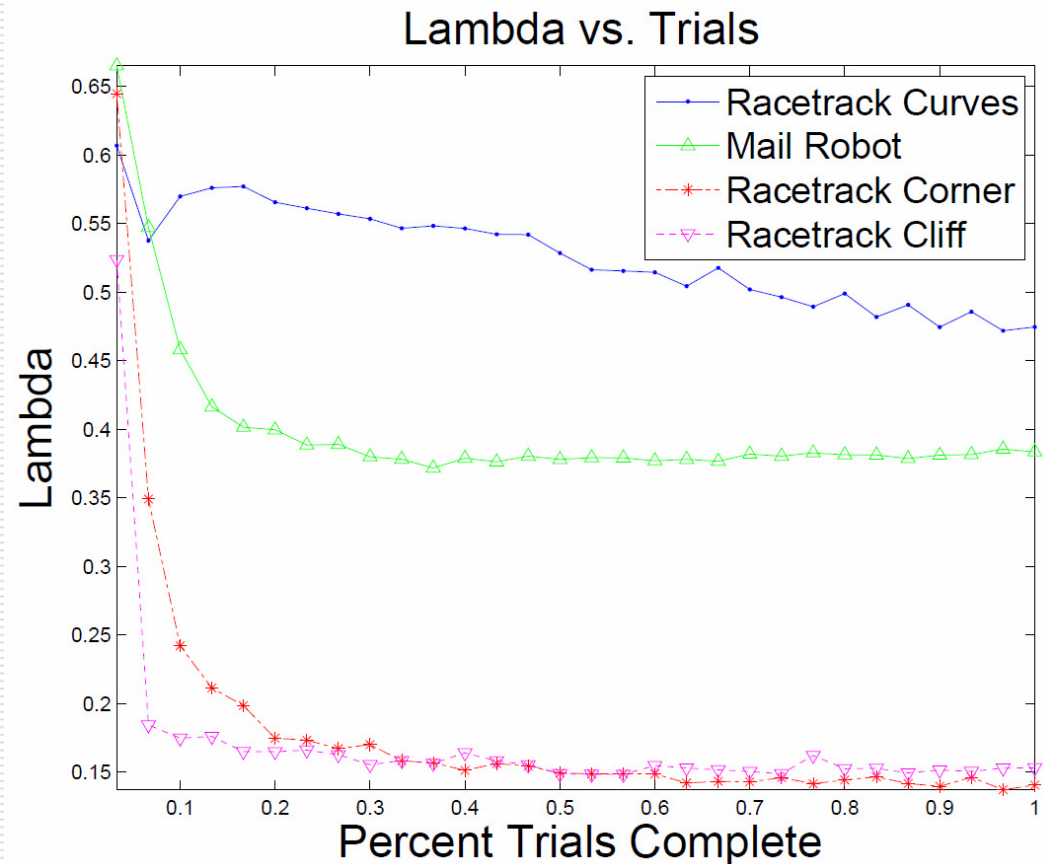
Results: Varying α



□ TD-BMA outperforms SARSA for all values of α .

Results: λ over time

- ❑ Bootstrapped model trusted more over time.
- ❑ Mail is partially observable, bootstrapped model helps disambiguate state.



Related Work

- Focus of other Bayesian RL approaches is on balancing exploration vs. Exploitation.
 - E.g. GPTD

 - Sutton and Singh propose three adaptive weighting schemes
 - First two restricted to acyclic.
 - Third requires maintaining an estimate of the model.
-

Conclusions and Future Work

- We derived a novel BMA approach to adapting λ in TD(λ).
 - We contributed the efficient Gaussian-based TD-BMA algorithm.
 - We showed that TD-BMA generally performs much better than SARSA for all fixed values of λ .
 - Future work:
 - Derive an online version of TD-BMA.
 - Use this for TD-BMA with function approximation.
-