Partially Observable RL with Memory Traces

Onno Eberhard¹² · Michael Muehlebach¹ · Claire Vernade²

¹Max Planck Institute for Intelligent Systems ²University of Tübingen

























- ▶ Memory is necessary in many partially observable environments
 - \rightarrow Memory: a compressed representation of the history of observations

- ▶ Memory is necessary in many partially observable environments
 - ightarrow Memory: a compressed representation of the history of observations
- ► Length-m window ("frame stacking"):

$$\operatorname{win}_{\mathfrak{m}}(\underbrace{y_{0}, y_{-1}, ...}_{\text{history } h}) \doteq (y_{0}, y_{-1}, ..., y_{-\mathfrak{m}+1})$$

- ▶ Memory is necessary in many partially observable environments
 - ightarrow Memory: a compressed representation of the history of observations
- ► Length-m window ("frame stacking"):

$$\operatorname{win}_{\mathfrak{m}}(\underbrace{y_{0}, y_{-1}, ...}_{\text{history } h}) \doteq (y_{0}, y_{-1}, ..., y_{-\mathfrak{m}+1})$$

▶ Our approach: the *memory trace* with forgetting factor $\lambda \in [0, 1)$:

$$z_{\lambda}(h) \doteq (1-\lambda) \sum_{k=0}^{|h|-1} \lambda^{k} y_{-k}$$

- ▶ Memory is necessary in many partially observable environments
 - ightarrow Memory: a compressed representation of the history of observations
- ► Length-m window ("frame stacking"):

$$\operatorname{win}_{\mathfrak{m}}(\underbrace{y_{0}, y_{-1}, \ldots}_{\text{history h}}) \doteq (y_{0}, y_{-1}, \ldots, y_{-\mathfrak{m}+1})$$

▶ Our approach: the *memory trace* with forgetting factor $\lambda \in [0, 1)$:

$$z_{\lambda}(h) \doteq (1-\lambda) \sum_{k=0}^{|h|-1} \lambda^{k} y_{-k}$$

 $\rightarrow \mbox{ Recursively computable: } z_{\lambda}(y_0,y_{-1},...) = \lambda z_{\lambda}(y_{-1},y_{-2},...) + (1-\lambda)y_0$























- \blacktriangleright We consider the problem of *policy evaluation* with offline data
 - \rightarrow The environment \mathcal{E} is a hidden Markov model (no explicit policy)
 - $\rightarrow~$ The observation space $\ensuremath{\mathfrak{Y}}$ is one-hot

- ▶ We consider the problem of *policy evaluation* with offline data
 - \rightarrow The environment \mathcal{E} is a hidden Markov model (no explicit policy)
 - $\rightarrow~$ The observation space $\ensuremath{\mathfrak{Y}}$ is one-hot
- ▶ Q: How much data do we need to accurately estimate the value function?

- \blacktriangleright We consider the problem of *policy evaluation* with offline data
 - \rightarrow The environment \mathcal{E} is a hidden Markov model (no explicit policy)
 - $\rightarrow~$ The observation space ${\mathfrak Y}$ is one-hot
- ▶ Q: How much data do we need to accurately estimate the value function?
- Given a function class $\mathfrak{F} \subset {\mathfrak{Y}^{\infty} \to [\underline{\nu}, \overline{\nu}]}$, find $f \in \mathfrak{F}$ that minimizes

$$\mathcal{R}_{\varepsilon}(f) \doteq \mathbb{E}_{\varepsilon} \left[\left\{ f(y_0, y_{-1}, \dots) - \sum_{t=0}^{\infty} \gamma^t r(y_{t+1}) \right\}^2 \right],$$

where $r: \mathcal{Y} \to [\underline{r}, \overline{r}]$ (rewards), $\gamma \in [0, 1)$ (discount), $\underline{\nu} \doteq \underline{r}/(1 - \gamma)$, and $\overline{\nu} \doteq \overline{r}/(1 - \gamma)$

- \blacktriangleright We consider the problem of *policy evaluation* with offline data
 - \rightarrow The environment \mathcal{E} is a hidden Markov model (no explicit policy)
 - $\rightarrow~$ The observation space ${\mathcal Y}$ is one-hot
- ▶ Q: How much data do we need to accurately estimate the value function?
- Given a function class $\mathfrak{F} \subset {\mathfrak{Y}^{\infty} \to [\underline{\nu}, \overline{\nu}]}$, find $f \in \mathfrak{F}$ that minimizes

$$\mathcal{R}_{\mathcal{E}}(f) \doteq \mathbb{E}_{\mathcal{E}}\left[\left\{f(y_0, y_{-1}, \dots) - \sum_{t=0}^{\infty} \gamma^t r(y_{t+1})\right\}^2\right],\$$

where $r: \mathcal{Y} \to [\underline{r}, \overline{r}]$ (rewards), $\gamma \in [0, 1)$ (discount), $\underline{\nu} \doteq \underline{r}/(1 - \gamma)$, and $\overline{\nu} \doteq \overline{r}/(1 - \gamma)$

• Window memory: $\mathcal{F}_{\mathfrak{m}} \doteq \{ \mathbf{f} \circ \operatorname{win}_{\mathfrak{m}} \mid \mathbf{f} : \mathcal{Y}^{\mathfrak{m}} \to [\underline{\nu}, \overline{\nu}] \}$

- \blacktriangleright We consider the problem of *policy evaluation* with offline data
 - \rightarrow The environment \mathcal{E} is a hidden Markov model (no explicit policy)
 - $\rightarrow~$ The observation space ${\mathcal Y}$ is one-hot
- ▶ Q: How much data do we need to accurately estimate the value function?
- Given a function class $\mathfrak{F} \subset {\mathfrak{Y}^{\infty} \to [\underline{\nu}, \overline{\nu}]}$, find $f \in \mathfrak{F}$ that minimizes

$$\mathcal{R}_{\varepsilon}(f) \doteq \mathbb{E}_{\varepsilon} \left[\left\{ f(y_0, y_{-1}, \dots) - \sum_{t=0}^{\infty} \gamma^t r(y_{t+1}) \right\}^2 \right],$$

where $r: \mathcal{Y} \to [\underline{r}, \overline{r}]$ (rewards), $\gamma \in [0, 1)$ (discount), $\underline{\nu} \doteq \underline{r}/(1 - \gamma)$, and $\overline{\nu} \doteq \overline{r}/(1 - \gamma)$

- Window memory: $\mathcal{F}_{\mathfrak{m}} \doteq \{ f \circ win_{\mathfrak{m}} \mid f : \mathcal{Y}^{\mathfrak{m}} \rightarrow [\underline{\nu}, \overline{\nu}] \}$
- Memory traces: $\mathfrak{F}_{\lambda} \doteq \{ f \circ z_{\lambda} \mid f : \mathfrak{Z}_{\lambda} \rightarrow [\underline{\nu}, \overline{\nu}] \}$, where $\mathfrak{Z}_{\lambda} \doteq \{ z_{\lambda}(h) \mid h \in \mathfrak{Y}^{\infty} \}$

The geometry of the *trace space* \mathcal{Z}_{λ}







The geometry of the *trace space* \mathcal{Z}_{λ}



The geometry of the *trace space* \mathcal{Z}_{λ}





Onno Eberhard · Partially Observable RL with Memory Traces



Onno Eberhard · Partially Observable RL with Memory Traces

The geometry of the *trace space* \mathcal{Z}_{λ}



The geometry of the *trace space* \mathcal{Z}_{λ}



0.8 8 $\lambda = 0.50$ $\mathcal{Y} = \{a, b, c\}$ Δ 0.6 ልን ልን የ $z_3)/\sqrt{6}$ 0.4 length 5 0.2 z_{2} 7 5 History 3 0.0 $(2z_1$ <u>ርእ ልእ</u> / -0.2-0.40 0.5 -0.5 $(z_2 - z_3)/\sqrt{2}$

Onno Eberhard · Partially Observable RL with Memory Traces

0.8 8 $\lambda = 0.50$ $\mathcal{Y} = \{a, b, c\}$ 0.6 $z_3)/\sqrt{6}$ ດາດາ 0.4 length 5 0.2 z_{2} 7 5 History 3 0.0 $(2z_1$ С b -0.2-0.40 0.5 -0.5 $(z_2 - z_3)/\sqrt{2}$

Onno Eberhard · Partially Observable RL with Memory Traces

Learning value functions

Theorem (Hoeffding bound)

Given a dataset \mathfrak{D} of a trajectories from an environment \mathcal{E} , a function class \mathcal{F} , and some c > 0, let \mathcal{F} be the smallest c-cover of \mathcal{F} and $f_{\mathfrak{m}} = \arg\min_{e \in \mathcal{F}} \sum_{\mathfrak{D}} (f(\mathfrak{h}) - \sum_{e \in \mathcal{F}} \mathcal{F}_{\mathfrak{m}}(\mathfrak{g}, \mathfrak{h}))$ Learning a good value estimate seems easier if the metric entropy $\mathcal{H}_{\epsilon}(\mathcal{F})$ of the function class \mathcal{F} is small. $\mathcal{R}_{\epsilon}(f_{\mathfrak{m}}) \leq \mathcal{R}_{\epsilon}(\mathcal{F}) + (\mathfrak{v} - \mathfrak{v})^{2} \sqrt{\frac{\mathcal{H}_{\epsilon}(\mathcal{F}) + \log 2}{2\mathfrak{n}}} + O(\epsilon).$

Let (X,p) be a metric space and I C X. The *e-covering number* N.(I) is the carThe metric entropy is a measure of "size" for the (infinite) function class. If the metric entropy of T is defined as H.(I) = log N.(I).

Learning value functions

Theorem (Hoeffding bound)

Given a dataset \mathcal{D} of n trajectories from an environment \mathcal{E} , a function class \mathcal{F} , and some $\varepsilon > 0$, let $\mathcal{F}^{\varepsilon}$ be the smallest ε -cover of \mathcal{F} and $f_n \doteq \arg\min_{f \in \mathcal{F}^{\varepsilon}} \sum_{\mathcal{D}} \{f(h) - \sum_{t=0}^{\infty} \gamma^t r(y_{t+1})\}^2$. Then, with probability at least $1 - \delta$,

$$\mathcal{R}_{\mathcal{E}}(f_{\mathfrak{n}}) \leq \mathcal{R}_{\mathcal{E}}(\mathcal{F}) + (\overline{\nu} - \underline{\nu})^2 \sqrt{\frac{\mathsf{H}_{\varepsilon}(\mathcal{F}) + \log \frac{2}{\delta}}{2\mathfrak{n}}} + \mathcal{O}(\varepsilon).$$

Let (X, ρ) be a metric space and $T \subset X$. The ϵ -covering number $N_{\epsilon}(T)$ is the cardinality of the smallest set $S \subset X$ such that for every $x \in T$, there exists a $y \in S$ with $\rho(x, y) \leqslant \epsilon$. The metric entropy of T is defined as $H_{\epsilon}(T) \doteq \log N_{\epsilon}(T)$.

► For windows of length m, we have $H_{\epsilon}(\mathcal{F}_m) \in \Theta(|\mathcal{Y}|^m)$

 $\rightarrow\,$ Exponential in m: long windows are expensive!

► For windows of length m, we have $H_{\epsilon}(\mathcal{F}_{\mathfrak{m}}) \in \Theta(|\mathcal{Y}|^{\mathfrak{m}})$

 $\rightarrow\,$ Exponential in m: long windows are expensive!

• Memory traces with forgetting factor $\lambda < \frac{1}{2}$ remember everything

 \rightarrow There is no compression: z_{λ} is invertible and therefore $H_{\epsilon}(\mathcal{F}_{\lambda}) = \infty$

- ► For windows of length m, we have $H_{\varepsilon}(\mathcal{F}_m) \in \Theta(|\mathcal{Y}|^m)$
 - $\rightarrow\,$ Exponential in m: long windows are expensive!
- Memory traces with forgetting factor $\lambda < \frac{1}{2}$ remember everything \rightarrow There is no compression: z_{λ} is invertible and therefore $H_{\epsilon}(\mathcal{F}_{\lambda}) = \infty$
- ▶ Need to "zoom in" to differentiate histories that only differ far in the past

- ► For windows of length m, we have $H_{\varepsilon}(\mathcal{F}_m) \in \Theta(|\mathcal{Y}|^m)$
 - $\rightarrow\,$ Exponential in m: long windows are expensive!
- Memory traces with forgetting factor $\lambda < \frac{1}{2}$ remember everything \rightarrow There is no compression: z_{λ} is invertible and therefore $H_{e}(\mathcal{F}_{\lambda}) = \infty$
- ▶ Need to "zoom in" to differentiate histories that only differ far in the past
- ▶ The "resolution" of a function class is given by its Lipschitz constant

- ► For windows of length m, we have $H_{\epsilon}(\mathcal{F}_{\mathfrak{m}}) \in \Theta(|\mathcal{Y}|^{\mathfrak{m}})$
 - $\rightarrow\,$ Exponential in m: long windows are expensive!
- Memory traces with forgetting factor $\lambda < \frac{1}{2}$ remember everything \rightarrow There is no compression: z_{λ} is invertible and therefore $H_{\varepsilon}(\mathcal{F}_{\lambda}) = \infty$
- ▶ Need to "zoom in" to differentiate histories that only differ far in the past
- ▶ The "resolution" of a function class is given by its Lipschitz constant
- We consider the class $\mathcal{F}_{\lambda,L} \doteq \{f \circ z_{\lambda} \mid f : \mathcal{Z}_{\lambda} \to [\underline{\nu}, \overline{\nu}], f \text{ is } L\text{-Lipschitz}\}$

- ► For windows of length m, we have $H_{\epsilon}(\mathcal{F}_{\mathfrak{m}}) \in \Theta(|\mathcal{Y}|^{\mathfrak{m}})$
 - $\rightarrow\,$ Exponential in m: long windows are expensive!
- Memory traces with forgetting factor $\lambda < \frac{1}{2}$ remember everything \rightarrow There is no compression: z_{λ} is invertible and therefore $H_{\varepsilon}(\mathcal{F}_{\lambda}) = \infty$
- ▶ Need to "zoom in" to differentiate histories that only differ far in the past
- ▶ The "resolution" of a function class is given by its Lipschitz constant
- ► We consider the class $\mathcal{F}_{\lambda,L} \doteq \{f \circ z_{\lambda} \mid f : \mathcal{Z}_{\lambda} \rightarrow [\underline{\nu}, \overline{\nu}], f \text{ is } L\text{-Lipschitz}\}$

Lemma (metric entropy of $\mathcal{F}_{\lambda,L}$)

$$\mathsf{H}_{\epsilon}(\mathcal{F}_{\lambda,L}) \in \mathcal{O}\big(L^{\min\{d_{\lambda},|\mathcal{Y}|-1\}}\big) \qquad \text{where} \qquad d_{\lambda} \doteq \frac{\log |\mathcal{Y}|}{\log(1/\lambda)}$$

Fast forgetting: $\lambda < \frac{1}{2}$

Theorem (window $\rightarrow ext{trace}$)

Let $\mathfrak{m} \in \mathbb{N}$ be a window length, $0 < \lambda < \frac{1}{2}$ a forgetting factor, and define

Windows are not more efficient than memory traces.

Then, for every $\epsilon > 0$ and every environment \mathcal{E} ,

 $\Re_{\mathcal{E}}(\mathcal{F}_{\lambda,L(\mathfrak{m})}) \leqslant \Re_{\mathcal{E}}(\mathcal{F}_{\mathfrak{m}}) \quad \text{and} \quad \mathsf{H}_{\epsilon}(\mathcal{F}_{\lambda,L(\mathfrak{m})}) \in \mathcal{O}(|\mathcal{Y}|^{\mathfrak{m}}) = \mathcal{O}(\mathsf{H}_{\epsilon}(\mathcal{F}_{\mathfrak{m}})).$

Fast forgetting:
$$\lambda < \frac{1}{2}$$

Theorem (window \rightarrow trace)

Let $m \in \mathbb{N}$ be a window length, $0 < \lambda < \frac{1}{2}$ a forgetting factor, and define

$$L(\mathfrak{m}) = \frac{\overline{\nu} - \underline{\nu}}{\sqrt{2}(1 - 2\lambda)\lambda^{\mathfrak{m} - 1}}$$

Then, for every $\epsilon > 0$ and every environment \mathcal{E} ,

 $\mathcal{R}_{\mathcal{E}}(\mathcal{F}_{\lambda,L(\mathfrak{m})}) \leqslant \mathcal{R}_{\mathcal{E}}(\mathcal{F}_{\mathfrak{m}}) \quad \text{and} \quad \mathsf{H}_{\varepsilon}(\mathcal{F}_{\lambda,L(\mathfrak{m})}) \in \mathcal{O}(|\mathcal{Y}|^{\mathfrak{m}}) = \mathcal{O}(\mathsf{H}_{\varepsilon}(\mathcal{F}_{\mathfrak{m}})).$

Fast forgetting: $\lambda < \frac{1}{2}$

Theorem (trace \rightarrow window)

Let $\lambda \in [0, 1)$ be a forgetting factor, L > 0 a Lipschitz constant, $\epsilon \in (0, L)$, and define

$\mathfrak{m}(\lambda, \mathbf{L}) = \left\lceil \frac{\log(\mathbf{L}/\epsilon)}{\log(1/\lambda)} \right\rceil.$

Then, Memory traces with $\lambda < \frac{1}{2}$ seem no more efficient than windows.

 $\mathcal{R}_{\mathcal{E}}ig(\mathcal{F}_{\mathfrak{m}(\lambda,\mathrm{L})}ig)\leqslant\mathcal{R}_{\mathcal{E}}(\mathcal{F}_{\lambda,\mathrm{L}})+\mathbb{O}(\epsilon) \quad ext{and} \quad \mathsf{H}_{\epsilon}ig(\mathcal{F}_{\mathfrak{m}(\lambda,\mathrm{L})}ig)\in\mathbb{O}(\mathsf{L}^{\mathsf{d}_{\lambda}}).$

If $\lambda < \frac{1}{2}$, then $d_{\lambda} < |\mathcal{Y}| - 1$.

Fast forgetting: $\lambda < \frac{1}{2}$

Theorem (trace \rightarrow window)

Let $\lambda \in [0, 1)$ be a forgetting factor, L > 0 a Lipschitz constant, $\epsilon \in (0, L)$, and define

$$\mathfrak{m}(\lambda, L) = \left| \frac{\log(L/\epsilon)}{\log(1/\lambda)} \right|.$$

Then, for every environment \mathcal{E} ,

$$\Re_{\mathcal{E}}(\mathfrak{F}_{\mathfrak{m}(\lambda,L)}) \leqslant \Re_{\mathcal{E}}(\mathfrak{F}_{\lambda,L}) + \mathfrak{O}(\epsilon) \quad \text{and} \quad \mathsf{H}_{\epsilon}(\mathfrak{F}_{\mathfrak{m}(\lambda,L)}) \in \mathfrak{O}(\mathsf{L}^{\mathsf{d}_{\lambda}}).$$

If $\lambda < \frac{1}{2}$, then $d_{\lambda} < |\mathcal{Y}| - 1$.

Slow forgetting: $\lambda \ge \frac{1}{2}$

Theorem (T-maze)

There exists a sequence (\mathcal{E}_k) of environments (with constant observation space \mathcal{Y}) with the property that, for every $\epsilon > 0$,

Memory traces $(\lambda \ge \frac{1}{2})$ can be significantly more efficient than windows.

In particular, the *T-maze* with corridor length k is such a sequence. In this case, the minima are attained at $m_k = k$, $\lambda_k = \frac{k-1}{k}$, and $L_k \leq \sqrt{2}ek$.

Slow forgetting: $\lambda \ge \frac{1}{2}$

Theorem (T-maze)

There exists a sequence (\mathcal{E}_k) of environments (with constant observation space \mathcal{Y}) with the property that, for every $\epsilon > 0$,

$$\min_{\mathfrak{m}\in\mathbb{N}} \{ \mathsf{H}_{\epsilon}(\mathcal{F}_{\mathfrak{m}}) \mid \mathcal{R}_{\mathcal{E}_{k}}(\mathcal{F}_{\mathfrak{m}}) = 0 \} \in \Omega(|\mathcal{Y}|^{\kappa}), \text{ and}$$
$$\min_{\lambda \in [0,1]} \min_{L \ge 0} \{ \mathsf{H}_{\epsilon}(\mathcal{F}_{\lambda,L}) \mid \mathcal{R}_{\mathcal{E}_{k}}(\mathcal{F}_{\lambda,L}) = 0 \} \in \mathcal{O}(k^{|\mathcal{Y}|-1}).$$

In particular, the *T*-maze with corridor length k is such a sequence. In this case, the minima are attained at $m_k = k$, $\lambda_k = \frac{k-1}{k}$, and $L_k \leq \sqrt{2}ek$.

• For $\lambda < \frac{1}{2}$, learning with windows and traces seems equivalent

- For $\lambda < \frac{1}{2}$, learning with windows and traces seems equivalent
- For $\lambda \ge \frac{1}{2}$, there exist environments where traces are much more efficient

Why can traces be more efficient?

- For $\lambda < \frac{1}{2}$, learning with windows and traces seems equivalent
- For $\lambda \ge \frac{1}{2}$, there exist environments where traces are much more efficient

Why can traces be more efficient?

 \blacktriangleright In the T-maze, most of the $| \mathfrak{Y} |^k$ histories are irrelevant \rightarrow allows for small L

- For $\lambda < \frac{1}{2}$, learning with windows and traces seems equivalent
- For $\lambda \ge \frac{1}{2}$, there exist environments where traces are much more efficient

Why can traces be more efficient?

- \blacktriangleright In the T-maze, most of the $|\mathfrak{Y}|^k$ histories are irrelevant \rightarrow allows for small L
- ▶ In other environments, memory traces can effectively smooth out noise



Deep RL with memory traces



Summary



- ▶ Memory is important in partially observable environments
- ► We analyze *memory traces*: exponential moving averages of past observations
- ▶ There is a close connection to *sliding window* memory
- ▶ We prove that memory traces are strictly more powerful than sliding windows
- ▶ Memory traces are an effective drop-in replacement for *frame stacking*
- ▶ Paper, code & more at onnoeberhard.com/memory-traces

Thank you for listening!

Bibliography

Bakker, Bram (2001). "Reinforcement Learning with Long Short-Term Memory". In: Advances in Neural Information Processing Systems. Vol. 14. LINK.