

Evaluación de políticas bajo ruido Markoviano mediante el algoritmo de Online Bootstrap Inference

Ana María Patrón Piñerez

10 de agosto de 2023

Introducción

- El objetivo es encontrar la función valor para una política dada, en casos donde el ambiente provee información insuficiente o sea de dimensiones grandes.
- La estrategia es aproximar linealmente la función valor [Sutton and Barto, 2018].
- Con aproximar linealmente una función nos referimos a aplicar Aproximación Lineal Estocástica (LSA). En el caso estándar, el ruido es una diferencia de martingalas; pero con nuestra estrategia, esto se incumple y ahora el ruido es Markoviano. [Ramprasad et al., 2022, Liang, 2010]

- Los métodos de Diferencias temporales (TD), para el caso on-policy, y de Gradiente de Diferencias Temporales (GTD), para el caso off-policy, son métodos clásicos de LSA para hacer evaluación de políticas bajo nuestra estrategia. Sin embargo, sus resultados son sesgados [Sutton and Barto, 2018, Maei, 2011].
- El algoritmo de Online Bootstrap plantea mejorar esto. Consiste en agregarle a los métodos clásicos un bootstrap estadístico, en el que para cada período de tiempo se realizan B iteraciones penalizadas del parámetro de LSA [Ramprasad et al., 2022].
- A partir de estas, se construyen intervalos de confianza para la función valor.

Contribución

- Reconstruir el contexto teórico en el que se construye el algoritmo de Online Bootstrap y compararlo con el escenario clásico.
- Robustecer el análisis para caso off-policy, usando Diferencias Temporales con corrección de Gradiente (TDC).
- Estudiar computacionalmente el desempeño del Online Bootstrap respecto a la política escogida.

Índice

- 1 Aproximación Lineal Estocástica (LSA)
 - Marco general
 - LSA con ruido Markoviano
 - El algoritmo de Online Bootstrap Inference
 - 2 Evaluación de políticas en Aprendizaje Reforzado (RL) usando LSA
 - Métodos clásicos
 - El algoritmo de Online Bootstrap Inference aplicado a RL
 - 3 Experimentos
- Bibliografía43

- 1 Aproximación Lineal Estocástica (LSA)
 - Marco general
 - LSA con ruido Markoviano
 - El algoritmo de Online Bootstrap Inference

 - 2 Evaluación de políticas en Aprendizaje Reforzado (RL) usando LSA
 - Métodos clásicos
 - El algoritmo de Online Bootstrap Inference aplicado a RL

 - 3 Experimentos
- Bibliografía43

Marco General

La Aproximación Lineal Estocástica (LSA) es un método iterativo para encontrar las raíces de una función lineal desconocida $f(x)$, a partir de unas observaciones $g(x_t; t+1)$, donde el ruido $t+1$ perturba las observaciones y se cumple que $f(x) = E[g(x; \cdot)]$

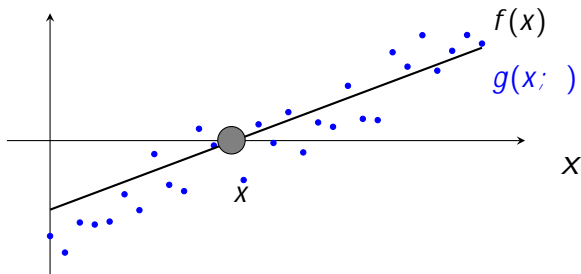


Figura: Representación gráfica de LSA.

Notacionalmente, tomamos

$$f(x_t) = \bar{A}x_t \quad \bar{b}; \quad g(x_t; t+1) = \tilde{A}_{t+1}x_t \quad \tilde{b}_{t+1}$$

donde

- $\bar{A}; \tilde{A}_t \in \mathbb{R}^{d \times d}; E[\tilde{A}_t] = \bar{A}$
- $\bar{b}; \tilde{b}_t; x_t \in \mathbb{R}^d; E[\tilde{b}_t] = \bar{b}$

Y buscamos encontrar x que soluciona la ecuación $\bar{A}x = \bar{b}$

Iteración de LSA (Robbins Monro)

Formalmente, la actualización está dada por:

$$x_{t+1} = x_t + \alpha_{t+1}[\tilde{A}_t x_t - \tilde{b}_t] \quad (1)$$

- x_t es el parámetro de LSA
- $\tilde{A}_t x_t - \tilde{b}_t$ es una observación con ruido de $\bar{A}x_t - \bar{b}$, donde el ruido $f(\tilde{A}_{t+1} - \bar{A})x_t - (\tilde{b}_{t+1} - \bar{b})g$ es una diferencia de martingalas
- α_{t+1} es un escalar positivo (step-size function en la literatura)

Teoría asintótica para LSA

Teorema 1 (Convergencia casi siempre del algoritmo de Robbins Monro)

- 1 Sea $e_{t+1} = (\tilde{A}_{t+1} - \bar{A})x_t - (\tilde{b}_{t+1} - \bar{b})$, donde $f e_t g_t = 0$ es una sucesión de diferencia de martingalas; es decir, $E[e_{t+1} | F_t] = 0$, donde F_t es el conjunto de información.
- 2 La sucesión $f \tilde{A}_{t+1} x_t - \tilde{b}_{t+1} g_t = 0 \in L^2$ con $E[|j \tilde{A}_{t+1} x_t - \tilde{b}_{t+1} j|^2 | F_t] \leq K(1 + |j x_t j|^2)$ c.s para algún $K > 0$
- 3 $\sum_{t=0}^{\infty} \frac{1}{t} = \infty$ y $\sum_{t=0}^{\infty} \frac{1}{t^2} < \infty$
- 4 La sucesión $f x_t g_t = 0$ esta acotada, $\sup_t |j x_t j|^2 < \infty$
 entonces $\lim_{t \rightarrow \infty} x_t = x$ c.s

Teoría asintótica para LSA

Teorema 2 (Normalidad asintótica de Robbins Monro)

Suponga que las condiciones 1-4 del teorema 1 valen. Además que

$$Q = \lim_{t \rightarrow \infty} \frac{1}{t} E[e_t e_t^T]$$

\bar{A} es una matriz Hurwitz, es decir que todos sus valores propios tienen parte real estrictamente negativa.

entonces

$$\rho_{\bar{t}}(x_t - x^*) \stackrel{d}{\rightarrow} N(0; \bar{A}^{-1} Q (\bar{A}^{-1})^T)$$

Corolario 3 (Polyak-Ruppert)

La convergencia casi siempre y la normalidad asintótica también valen para la iteración del estimador promediado $\bar{x}_t = \frac{1}{t} \sum_{i=1}^t x_i$.

Construcción de LSA con ruido Markoviano

Definición 4 (Cadena de Markov)

Una cadena de Markov es un proceso estocástico discreto $\{X_t\}_{t=0}^{\infty}$ tal que

$$P[X_{t+1}|X_t] = P[X_{t+1}|X_0; \dots; X_t]$$

Su kernel de transición P es una matriz con las probabilidades de transición de un evento a otro.

Teorema 5

Una cadena ergódica de Markov tiene distribución límite, que es única y es igual a su distribución estacionaria.

Definición 6 (LSA con ruido markoviano)

Sea $f_t, \tilde{A}(X_t); \tilde{b}(X_t)g_{t-1}$ una secuencia de observaciones, donde f_t, X_t, g_t es una cadena ergódica de Markov con espacio de estados X y distribución estacionaria π . Además $\tilde{A}: X \rightarrow \mathbb{R}^d$; $\tilde{b}: X \rightarrow \mathbb{R}^d$ y $E[\tilde{A}(X_{t+1})] = \bar{A}$; $E[\tilde{b}(X_{t+1})] = \bar{b}$. Entonces la actualización dada por

$$x_{t+1} = x_t + \alpha_{t+1}(\tilde{A}(X_{t+1})x_t - \tilde{b}(X_{t+1})) \quad (2)$$

corresponde al algoritmo de LSA con ruido Markoviano, donde $x \in \mathbb{R}^d$ es el parámetro de LSA y f_t, g_{t-1} es una sucesión de pasos decrecientes, dada por $\alpha_t = \frac{\alpha_0}{t}$ para algún $\alpha_0 > 0$; $\alpha \in (\frac{1}{2}, 1)$

Teoría asintótica para LSA con ruido Markoviano

- Ahora, $t_{+1} = (\tilde{A}(X_{t+1}) - \bar{A})x_t + (\tilde{b}(X_{t+1}) - \bar{b})$ y no se puede garantizar que sea una diferencia de martingalas
- Pero podemos escribir $t = e_t + \alpha t + \beta t$, donde e_t es la diferencia de martingalas de LSA en el caso general, t y t son términos residuales decrecientes [Liang, 2010].
- Si para la definición 6, además se tiene que \bar{A} es de rango completo, Hurwitz y que existen $A_{max}; b_{max}$ tales que $\sup_{x \in \mathcal{X}} \|\tilde{A}(x)\|_F \leq A_{max}; \sup_{x \in \mathcal{X}} \|\tilde{b}(x)\|_2 \leq b_{max}$, entonces:

Proposición 1

La convergencia casi siempre y la normalidad asintótica valen para la estimación de LSA con ruido Markoviano, en particular,

$$\sqrt{t}(\bar{\theta}_t - \theta^*) \xrightarrow{d} N(0; \bar{A}^{-1} Q (\bar{A}^{-1})^T) \text{ donde } Q = \lim_{t \rightarrow \infty} \frac{1}{t} E[e_t e_t^T]$$

- Aún así, en la práctica no podemos obtener valores para cada componente de $\bar{A}^{-1} Q (\bar{A}^{-1})^T$ es desconocido.
- Una alternativa es usar el algoritmo de Online Bootstrap Inference. [Ramprasad et al., 2022]

Su aporte es adicionar una iteración perturbada paralela a de la iteración promedio de LSA:

$$\begin{aligned} \hat{x}_{t+1}^b &= \hat{x}^b + W_{t+1}^b (A(\hat{x}_{t+1}) \hat{x}^b - b(\hat{x}_{t+1})) \\ \bar{x}_{t+1}^b &= \frac{1}{t+1} \sum_{i=1}^{t+1} \hat{x}_i^b \end{aligned} \quad (3)$$

donde W_t^b son variables aleatorias acotadas i.i.d con media y varianza 1.

¿La razón? se puede mostrar que

$\bar{x}_t(\hat{x}_t - \bar{x}_t) \stackrel{p}{\rightarrow} N(0; \bar{A}^{-1} Q (\bar{A}^{-1})^T); t \rightarrow \infty$ (Teo. 4.2 [Ramprasad et al., 2022]) ... <Es la misma de $\bar{x}_t - x$!

Podemos aproximar la distribución de $\bar{x}_t(\hat{x}_t - \bar{x}_t)$ con B muestras bootstraps de $\bar{x}_t - \bar{x}_t$ para cada t , es decir (3).

Algoritmo 1 Online Bootstrap Inference

Input: Numero de muestras bootstrap B , α_0 ; $\beta \in (0, 1)$ y
 estimadores iniciales $\hat{x}_0^b = \alpha_0^b$; $b = 1; \dots; B$ linesize=

```

for t = 0; 1; 2; ::: do
    Observar  $A(X_{t+1})$ ;  $b(X_{t+1})$ 
    Computar  $\beta_{t+1} = \beta(t+1)$ 
    Actualizar  $x_{t+1} = \beta x_t + \beta_{t+1} (A(X_{t+1}) - x_t) + b(X_{t+1})$ 
    Actualizar  $\bar{x}_{t+1} = \frac{1}{t+1} (t x_t + x_{t+1})$ 
    for b = 1; 2; :::; B do
        Actualizar  $\hat{x}_{t+1}^b = \beta x_t^b + \beta_{t+1} W_{t+1}^b (A(X_{t+1}) - x_t^b) + b(X_{t+1})$ 
        Actualizar  $\bar{x}_{t+1}^b = \frac{1}{t+1} (t \bar{x}_t^b + \hat{x}_{t+1}^b)$ 
    end for
end for
    
```

Output Estimadores bootstrap \hat{x}_{t+1}^b $g_{b=1}^B$

1 Aproximación Lineal Estocástica (LSA)

Marco general

LSA con ruido Markoviano

El algoritmo de Online Bootstrap Inference

2 Evaluación de políticas en Aprendizaje Reforzado (RL) usando LSA

Métodos clásicos

El algoritmo de Online Bootstrap Inference aplicado a RL

3 Experimentos

Bibliografía43

Figura: Figura tomada de
[Sutton and Barto, 2018]

Durante cada $t = 1; 2; \dots$ un agente se encuentra en un estado s_t .

Toma una acción a_t , de las disponibles para ese s_t .

Dependiendo de la configuración del ambiente, recibe una recompensa r_{t+1} y llega a un estado s_{t+1} .

Un Proceso de Decisión de Markov (MDP) es formalizar el ambiente de RL. Se denota $M = (S; A; P; R; \gamma)$ donde $S; A$ son los espacios de estados y acciones, respectivamente; P es el kernel de transición y R la función de recompensas.

Una política es un vector de probabilidades que para cada s_t , determina las probabilidades de tomar cada una de las acciones a_t disponibles.

Un MDP junto con una política inducen un Proceso de Recompensa de Markov (MRP), que se denota $M^\pi = (M; \pi)$.

La función valor asociada a π es la suma descontada de recompensas esperadas partiendo de s_0 y siguiendo π :

$$\begin{aligned} V^\pi(s) &= E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s; \pi\right]; \quad s \in S \\ &= R(s) + \gamma P^\pi(s) V^\pi(s) \\ &:= TV^\pi(s) \quad (\text{Ecuación de Bellman}) \end{aligned}$$

donde $P^\pi(s); R^\pi(s)$ son matrices calculables. Como $TV^\pi(s)$ es una contracción, por teorema de punto fijo de Banach $V^\pi(s)$ es la solución única. As:

$$V^\pi(s) = (I - \gamma P^\pi(s))^{-1} R^\pi(s) : \quad (4)$$

En RL, es común que haya incertidumbre, eso implica desconoce P .

Aun si hubiese certidumbre, las dimensiones del ambiente pueden ser grandes y así las dimensiones P de R
) <No se puede usar la ecuación de Bellman!

Una alternativa es aproximar linealmente la función valor como $V(s) = \phi(s)^T w$, donde $w \in \mathbb{R}^d$ es el parámetro de LSA y $\phi(s)^T = [\phi_1; \dots; \phi_d]$ con $\phi_i : S \rightarrow \mathbb{R}$ features (o funciones base)

Cuando se aplica LSA a RL con esa aproximación $V(s)$ el ruido de LSA es Markoviano. Aproximación: $X_t = f(s_t; r_{t+1}; s_{t+1}) g_t$.

I. Diferencias Temporales (TD)-On policy

La actualización de la función valor está dada por:

$$V_{t+1}(s_t) = V_t(s_t) + \alpha [V_t(s_t) - V_t(s_{t+1}) + r_{t+1}]: \quad (5)$$

Aplicando nuestra aproximación lineal de la función valor (5) es equivalente a:

$$V_{t+1} = V_t + \alpha [\underbrace{A(X_t)}_{(s_t)(s_t) \dots (s_{t+1})^T} \underbrace{b(X_t)}_{r_{t+1}(s_t)}]: \quad (6)$$

Que es la iteración de LSA con:

$$A(X_t) = (s_t)(s_t) \dots (s_{t+1})^T \text{ y } b(X_t) = r_{t+1}(s_t).$$

$$\bar{A} = E[A(X_t)] = \sum_{g \in \mathcal{G}} \pi(g) (I - P_g)^T \text{ y}$$

$$\bar{b} = E[b(X_t)] = \sum_{g \in \mathcal{G}} \pi(g) r_g, \text{ donde } \pi \text{ es la distribución estacionaria de } X_t \text{ g bajo la política } g.$$

Ahora:

En el caso on-policy, las observaciones se generan con la misma política a evaluar.

Pero en el caso off-policy, las observaciones se generan con una política b (behavior policy) distinta a la a evaluar (target policy).

Los problemas:

TD no es un gradiente inestabilidad.

Hay una desconexión, b puede seleccionar acciones que nunca tomara) inestabilidad.

Solución:

Metodos de Gradiente de Diferencia Temporales (GTD):
GTD1, GTD2 y TDC [Maei, 2011].

II. Gradiente de Diferencias Temporales (GTD)-O policy.

A. Derivación: usando gradiente descentente,

GTD1 minimiza la norma esperada de la actualización de TD:

$$NEU(\theta) = E [\delta_t(\theta) (\mathbf{s}_t)]^T E [\delta_t(\theta) (\mathbf{s}_t)]$$

donde $\delta_t(\theta) = r_{t+1} + \gamma V_t(\mathbf{s}_{t+1}) - V_t(\mathbf{s}_t)$ es el error de TD.

GTD2 y TDC minimizan el error cuadrático medio de Bellman proyectado:

$$MSBPE(\theta) = \sum_{jj} V \quad TV \quad jj^2$$

donde P es el operador de proyección, V matriz cuyos elementos en la diagonal son las entradas de

Al calcular el gradiente, se obtienen productos de valores esperados no necesariamente independientes. Se incluye una actualización auxiliar enlazada.

Se incluye un ratio de Important Sampling $w_t = \frac{p(a_t | s_t)}{q(a_t | s_t)}$ para corregir la discrepancia entre p y q :

Formalmente, tome

$$A_t = \sum_{a \in \mathcal{A}} w_t(a | s_t) (Q(s_t, a) - Q(s_t)) \quad b_t = \sum_{a \in \mathcal{A}} w_t(a | s_t) r_{t+1}(s_t)$$

como estimadores insesgados $A_t = \sum_{a \in \mathcal{A}} w_t(a | s_t) (Q(s_t, a) - Q(s_t))$ y $b_t = \sum_{a \in \mathcal{A}} w_t(a | s_t) r_{t+1}(s_t)$, respectivamente; entonces los algoritmos GTD se escriben como:

GTD1:

$$\begin{aligned} y_{t+1} &= y_t + \alpha_{t+1} [b_t + A_t y_t - y_t] \\ \theta_{t+1} &= \theta_t + \alpha_{t+1} A_t^T y_t \end{aligned} \quad (7)$$

GTD2:

$$\begin{aligned} y_{t+1} &= y_t + \alpha_{t+1} [b_t + A_t y_t - (\mathbf{s}_t)^T y_t - (\mathbf{s}_t)] \\ \theta_{t+1} &= \theta_t + \alpha_{t+1} A_t^T y_t \end{aligned} \quad (8)$$

TDC:

$$\begin{aligned} y_{t+1} &= y_t + \alpha_{t+1} [b_t + A_t y_t - (\mathbf{s}_t)^T y_t - (\mathbf{s}_t)] \\ \theta_{t+1} &= \theta_t + \alpha_{t+1} [b_t + A_t y_t - \theta_t (\mathbf{s}_{t+1}) (\mathbf{s}_t)^T y_t] \end{aligned} \quad (9)$$

B.Compactificando como problema de LSA:

GTD1 y GTD2:

$$\tilde{A}_t = \begin{pmatrix} 0 & A_t^T \\ A_t & M_t \end{pmatrix} \quad \tilde{b}_t = \begin{pmatrix} 0 \\ b_t \end{pmatrix} \quad \tilde{y}_t = \begin{pmatrix} r_t \\ y_t \end{pmatrix} \quad (10)$$

donde $M_t = I_d$ para GTD1 y $M_t = (s_t)^\top (s_t)$ para GTD2.

TDC:

$$\tilde{A}_t = \begin{pmatrix} A_t & (M_t + A_t)^\top \\ \alpha A_t & \alpha M_t \end{pmatrix} \quad \tilde{b}_t = \begin{pmatrix} b_t \\ \alpha b_t \end{pmatrix} \quad \tilde{y}_t = \begin{pmatrix} r_t \\ y_t \end{pmatrix} \quad (11)$$

donde $M_t = (s_t)^\top (s_t)$, α proxy empírica de $\lambda > 0$, donde λ es tal que es mayor al negativo del mínimo valor propio de la matriz $\begin{pmatrix} A & (A+A^\top) \\ \alpha A & \alpha(A+A^\top) \end{pmatrix}$.

Para un estado particular s ; $V_t^-(s) = \bar{V}_t^-(s)$, donde \bar{V}_t^- se obtiene de:

Algoritmo 2 Método clásico

Input: v_0 ; $\gamma \in (0, 1)$ y estimador inicial \bar{v}_0

for $t = 0; 1; 2; \dots$ do

Observar $A(X_{t+1})$; $b(X_{t+1})$

Computar $\bar{v}_{t+1} = v_0(t+1)$

Actualizar $\bar{v}_{t+1} = \gamma \bar{v}_t + (1-\gamma)(A(X_{t+1}) - b(X_{t+1}))$

Actualizar $\bar{v}_{t+1} = \frac{1}{1+t}(\bar{v}_t + \bar{v}_{t+1})$

end for

Output Estimador \bar{v}_{t+1}

1. Inicialmente, de la parte I, del algoritmo Online Bootstrap obtenamos $\bar{\lambda}_{t+1}^b g_{b=1}^B$
2. Ajustando el algoritmo a nuestro escenario de RL, para cada tenemos que $V_t^-(s) = T(s)^-_t$. As, ahora construimos

$$V_t^-(\cdot) = \int_{s \in S} V_t^-(s) \mu(ds)$$

3. De nimos los bootstraps de valores estimados como

$$\left(V_{\bar{\lambda}_t^{(b)}} \right)_{b=1}^B = V_t^- \quad (12)$$

4. Construimos intervalos de confianza con ancha.

$IC_{\text{cuantil}}(V_t(\cdot)) = [V_t(\cdot) + q_{\frac{\alpha}{2}}; V_t(\cdot) + q_{1-\frac{\alpha}{2}}]$,
 donde q es el q -ésimo cuantil de la distribución (12).

$IC_{\text{SE}}(V_t(\cdot)) = [V_t(\cdot) + z_{\frac{\alpha}{2}} \hat{\sigma}; V_t(\cdot) + z_{1-\frac{\alpha}{2}} \hat{\sigma}]$,
 donde z hace referencia a la distribución normal y
 $\hat{\sigma} = \frac{s^2}{B}$; s^2 es la varianza de (12).

1 Aproximación Lineal Estocástica (LSA)

Marco general

LSA con ruido Markoviano

El algoritmo de Online Bootstrap Inference

2 Evaluación de políticas en Aprendizaje

Reforzado (RL) usando LSA

Métodos clásicos

El algoritmo de Online Bootstrap Inference aplicado a RL

3 Experimentos

Bibliografía43

Evaluamos el desempeño del Online Bootstrap respecto al método clásico y su sensibilidad respecto a la política evaluada.

El código que se emplea es principalmente construido por [Ramprasad et al., 2022]. Adaptamos este para el caso con TDC y para variaciones en la política.

Los intervalos de confianza son al 95%; $n = 200$,
 $\alpha = \frac{2}{3}$; $\theta = \theta$. Además, $W_t \sim U[1 - \rho^{\frac{1}{3}}; 1 + \rho^{\frac{1}{3}}]$.

Analizamos un caso on-policy y otro off-policy.

Se usa el ambiente FrozenLake de OpenAI gym para generar el MDP.

Hay 64 estados, cada estado es una cuadrícula de una grilla 8×8 ; el espacio de acciones es

$A = \{f, d, r, a\}$ izquierda, derecha, arriba, abajo.

El objetivo es llegar al último estado (casilla con coordenadas $(8, 8)$).

La recompensa es 1 en el último estado y 0 de lo contrario.

La política, con la que se generan las observaciones, es una Q-política greedy con $\epsilon = 0.2$; y entre mayor sea, hay más aleatoriedad.

Función valor
verdadera

0;011(2 IC
siempre)

TD se acerca a
mayor número
de iteraciones

Hipotesis: Mayor distancia
entre el kernel de transición de
 X_t bajo aleatoria y su
distribución estacionaria se
relaciona con convergencia mas
lenta

= 0;3

= 0;5

Simulador del manejo de la sepsis (Oberst y Sontag, 2019).

Este simula 4 signos del paciente: frecuencia cardiaca, presión sanguínea, concentración de oxígeno y niveles de glucosa.

Hay 720 estados. Hay 8 acciones, que son las combinaciones de tratamientos posibles.

La recompensa es 1 si el paciente es dado de alta, -1 si muere y 0 en cualquier otro caso.

es óptima bajo Q-learning y Q^b es una Q-política greedy con $\gamma = 0.95$.

El Online Bootstrap tiene una motivación clara. Esta es la necesidad de evaluar políticas usando una estrategia conveniente y fácilmente implementable (en nuestro caso es aproximar linealmente la función valor), pero penalizando el ruido subyacente que distorsiona los resultados.

El Online Bootstrap permite garantizar que la función valor verdadera está dentro de un intervalo siempre. Aunque, TD tiene un comportamiento aceptable, su estimación no siempre es precisa.

El comportamiento de los métodos GTD es muy parecido. El que mejor se comporta es GTD2, le sigue TDC y por último GTD1. Es probable que estos efectos se deban a la función objetivo que están minimizando.

El algoritmo de Online Bootstrap es sensible ante cambios en la política evaluada. Aunque la convergencia se da, esta es más lenta y el tiempo computacional requerido es mayor.

- [Aré, 2017] [Aré, A. \(2017\)](#).
Stochastic approximation and martingale methods.
- [Bach et al., 2016] [Bach, F., Liu, Y., and Li, R. \(2016\)](#).
Statistical machine learning and convex optimization.
[Département d'Informatique de l'ENS \(DI ENS\)](#).
- [Bercu, 2019a] [Bercu, B. \(2019a\)](#).
Asymptotic behavior of stochastic algorithms with statistical applications. part i.
- [Bercu, 2019b] [Bercu, B. \(2019b\)](#).
Asymptotic behavior of stochastic algorithms with statistical applications. part ii.
- [Borkar, 2006] [Borkar, V. S. \(2006\)](#).
Stochastic approximation with 'controlled markov noise'.
[Systems and Control Letters, 55\(2\):139{145](#).
- [Borkar, 2008] [Borkar, V. S. \(2008\)](#).
Stochastic approximation a dynamical systems viewpoint.
[Cambridge University Press](#).
- [Haskell, 2011] [Haskell, W. B. \(2011\)](#).
Introduction to dynamic programming.
[University of Alberta](#).
- [Karmakar, 2020] [Karmakar, P. \(2020\)](#).
Stochastic approximation with markov noise: Analysis and applications in reinforcement learning.
- [Kushner and Yin, 2003] [Kushner, H. and Yin, G. \(2003\)](#).
Stochastic approximation and recursive algorithms and applications.
[Springer](#).

- [Levin et al., 2017] Levin, D., Peres, Y., and Wilmer, E. L. (2017).
Markov chains and mixing times.
American Mathematical Society.
- [Liang, 2010] Liang, F. (2010).
Trajectory averaging for stochastic approximation mcmc algorithms.
The Annals of Statistics.
- [Maei, 2011] Maei, H. R. (2011).
Gradient temporal-difference learning algorithms.
University of Alberta.
- [Oberst and Sontag, 2019] Oberst, M. and Sontag, D. (2019).
Counterfactual off-policy evaluation with gumbel-max structural causal models.
Proceedings of the 36th International Conference on Machine Learning.
- [Ramprasad et al., 2022] Ramprasad, P., Li, Y., Yang, Z., Wang, Z., Sun, W., and Cheng, G. (2022).
Online bootstrap inference for policy evaluation in reinforcement learning.
Journal of the American Statistical Association.
- [Robbins and Monroe, 1951] Robbins, H. and Monroe, S. (1951).
A stochastic approximation method.
The Annals of Mathematical Statistics.
- [Sutton and Barto, 2018] Sutton, R. and Barto, A. (2018).
Reinforcement learning: An introduction.
The MIT Press.
- [Xu et al., 2020] Xu, T., Wang, Z., Zhou, Y., and Liang, Y. (2020).
Reanalysis of variance reduced temporal difference learning.
International Conference on Learning Representations.