

APRENDIZAJE REFORZADO

CLASE 5

Julían Martínez

ESQUEMA "GENERAL"

$$v^{l+1}(S_t^{l+1}) = v^l(S_t^{l+1}) + \alpha \left[A(\text{episodo } l+1, v^l) - v^l(S_t^{l+1}) \right]$$

MC:

- NO USA LA MARKOVIANEIDAD
- PRECISA TODO EL EPISODIO PARA ACTUALIZAR LA FUNCIÓN DE VALOR.
- VARIANZA.

TD:

- USA EL SUPUESTO DE MARKOVIANEIDAD.
- ACTUALIZA LA FUNCIÓN EN CADA PASO DEL SAMPLING.
- PARA CADA ACTUALIZACIÓN UTILIZA LA ESTIMACIÓN PREVIA DE LA FUNCIÓN DE VALOR (SESGO).

GLIE CONTROL MONTE CARLO


$$\lim_{k \rightarrow \infty} N_k(s, a) = \infty$$

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = \delta_{\underset{a'}{\operatorname{argmax}}} q^+(s, a')$$

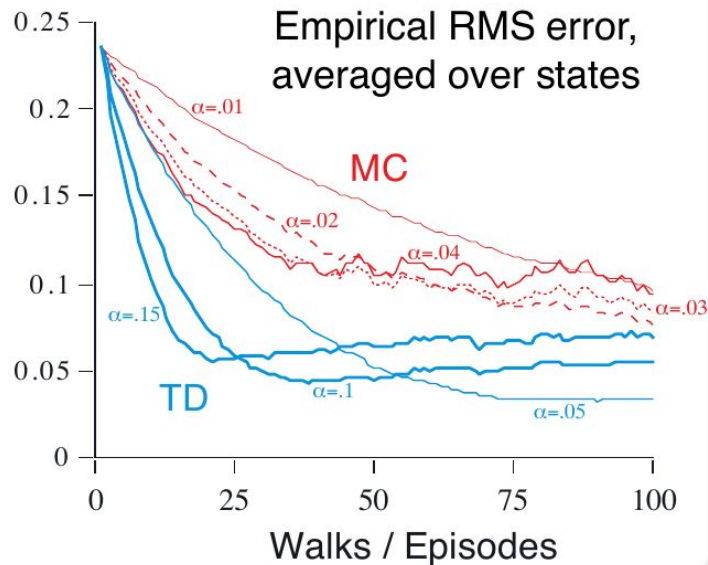
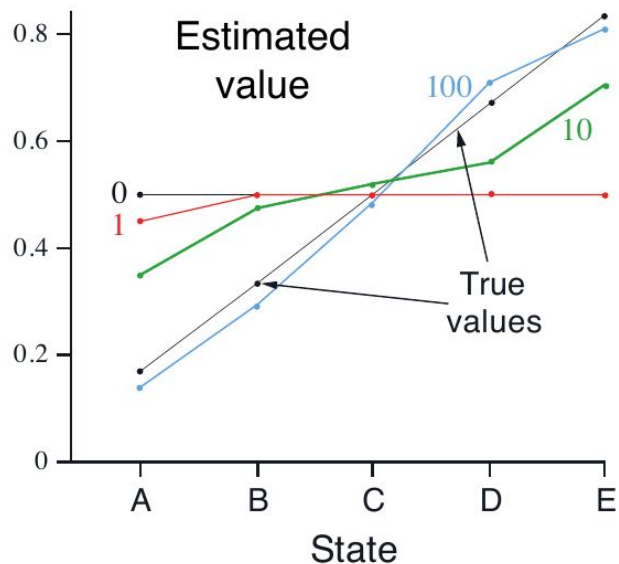
$$\varepsilon^l = \frac{1}{l}$$

PASA EN LAS MEJORES FAMILIAS

- The policy converges on a greedy policy,

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = \mathbf{1}(a = \operatorname{argmax}_{a' \in \mathcal{A}} Q_k(s, a'))$$


RANDOM WALK EXAMPLE 6.2 (MARKOV REWARD PROCESS)



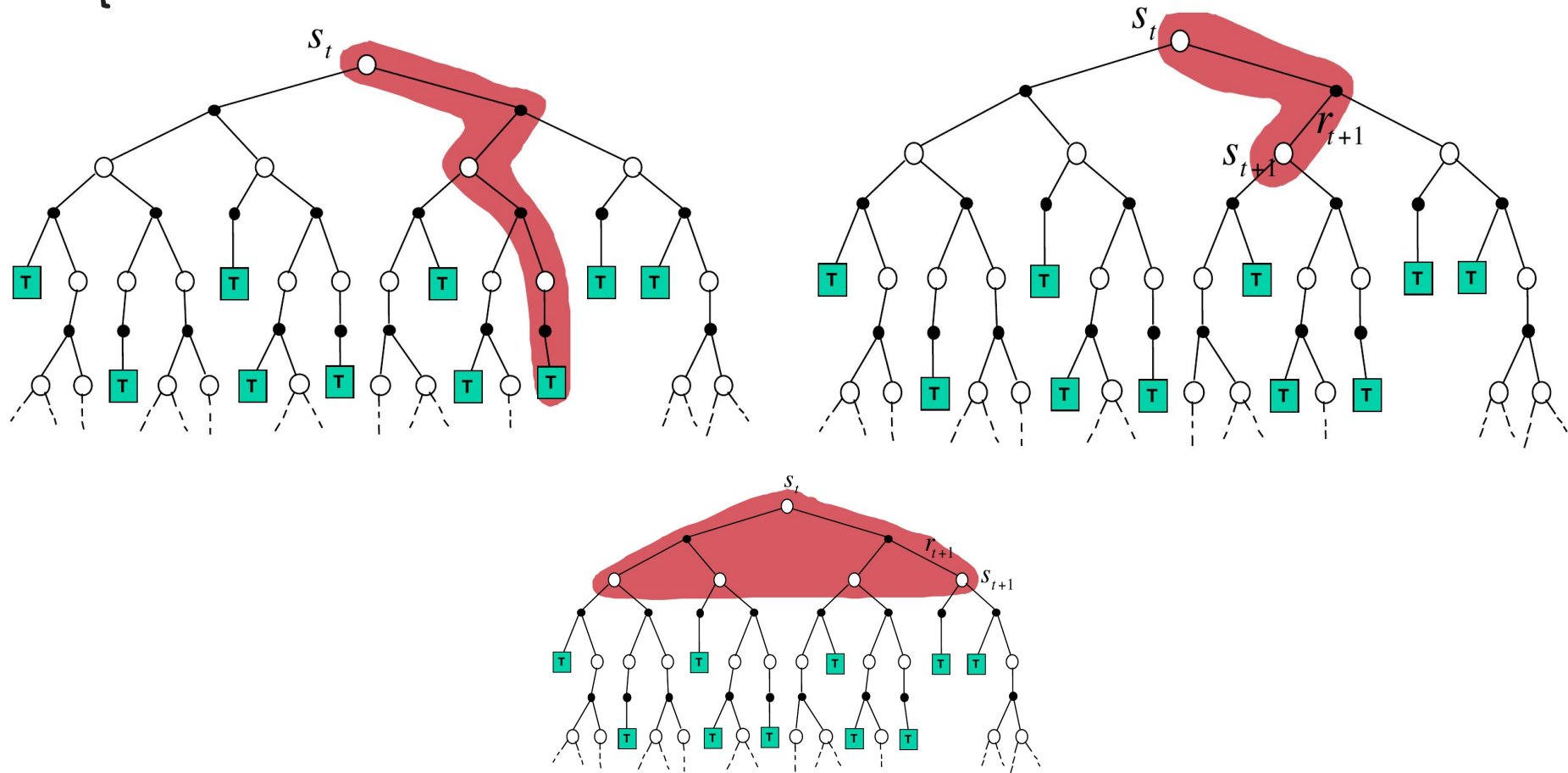
$$\sqrt{\sum_{r=1}^{100} \left\{ \frac{\sum_s [\hat{v}_r^l(s) - v(s)]^2}{|S|} \right\} \cdot \frac{1}{100}}$$

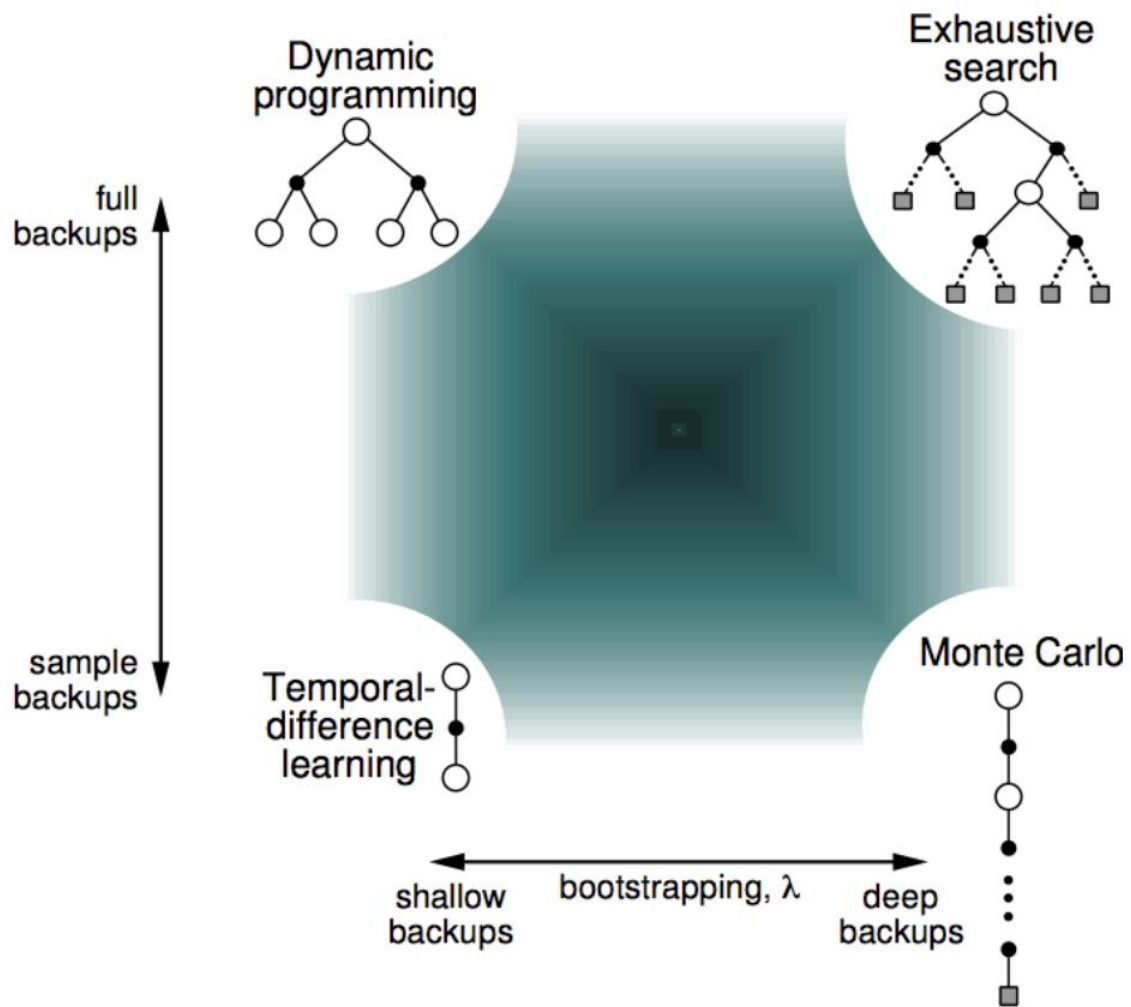
RELACIÓN ENTRE LOS ERRORES

$$G_t^{k+1} - v^k(s_t) = \dots = \sum_{j=t}^{T-1} \gamma^{j-t} \delta_j^{\text{TD}}$$

$$\begin{aligned}
G_t - V(S_t) &= R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1}) && \text{(from (3.9))} \\
&= \delta_t + \gamma(G_{t+1} - V(S_{t+1})) \\
&= \delta_t + \gamma\delta_{t+1} + \gamma^2(G_{t+2} - V(S_{t+2})) \\
&= \delta_t + \gamma\delta_{t+1} + \gamma^2\delta_{t+2} + \cdots + \gamma^{T-t-1}\delta_{T-1} + \gamma^{T-t}(G_T - V(S_T)) \\
&= \delta_t + \gamma\delta_{t+1} + \gamma^2\delta_{t+2} + \cdots + \gamma^{T-t-1}\delta_{T-1} + \gamma^{T-t}(0 - 0) \\
&= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k. && (6.6)
\end{aligned}$$

ESQUEMA DE BACKUPS





ENTRE MC Y TD

1-step TD
and TD(0)



2-step TD



3-step TD



...

n-step TD



...

∞ -step TD
and Monte Carlo



$$G_{t:t+n} := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n v(s_{t+n})$$

N-STEP TD

$$G_t^{(n)} \quad G_{t:t+n} := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n v(S_{t+n})$$

$$v^{t+n}(S_t) = v^{t+n-1}(S_t) + \alpha [G_{t:t+n} - v^{t+n-1}(S_t)]$$

ERROR REDUCTION PROPERTY

$$\begin{aligned} \max_s |E_{\pi}[G_{t:t+n} | S_t = s] - v_{\pi}(s)| \\ \leq \gamma^n \max_s |V_{t+n-1}(s) - v_{\pi}(s)| \end{aligned}$$

TD(λ) - FORMA "INGENIOSA" DE MEZCLAR LOS N-STEP

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

$$G_t^{(n)} = \sum_{k=1}^n \gamma^{k-1} R_{t+k} + \gamma^n \hat{v}(S_{t+n})$$

OBSERVACIÓN "ALGORÍTMICA"

$$\sum_{n=1}^{\infty} \left\{ \sum_{k=1}^n \lambda^{n-1} \cdot \gamma^{k-1} \cdot R_{t+k} + \lambda^{n-1} \hat{v}(S_{t+n}) \right\} (1-\lambda)$$

$$= \sum_{k=1}^{\infty} \sum_{n=1}^k (1-\lambda) \lambda^{n-1} \gamma^{k-1} R_{t+k} + (1-\lambda) \lambda^n \hat{v}(S_{t+n})$$

ESOS PESOS SON LA GEOMÉTRICA!

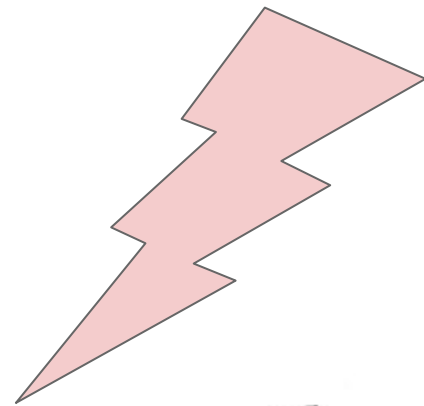


$$G_t^\lambda = (1-\lambda) \cdot \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

$$G_t^\lambda := E_\lambda[G_t^{(M)}]$$

$$M \sim \text{Geo}(1-\lambda)$$

QUE ELEGANCIA LA DE FRANCIA...



$$G_t = (1-\gamma)^{-1} E_{\gamma}[R_{\tau+t}]$$

$$\tau \sim \text{Geo}(1-\gamma)$$

REINTERPRETANDO LA FUNCIÓN DE VALOR

$$v(s) := \mathbb{E}_{\pi} [G_t | S_t = s]$$

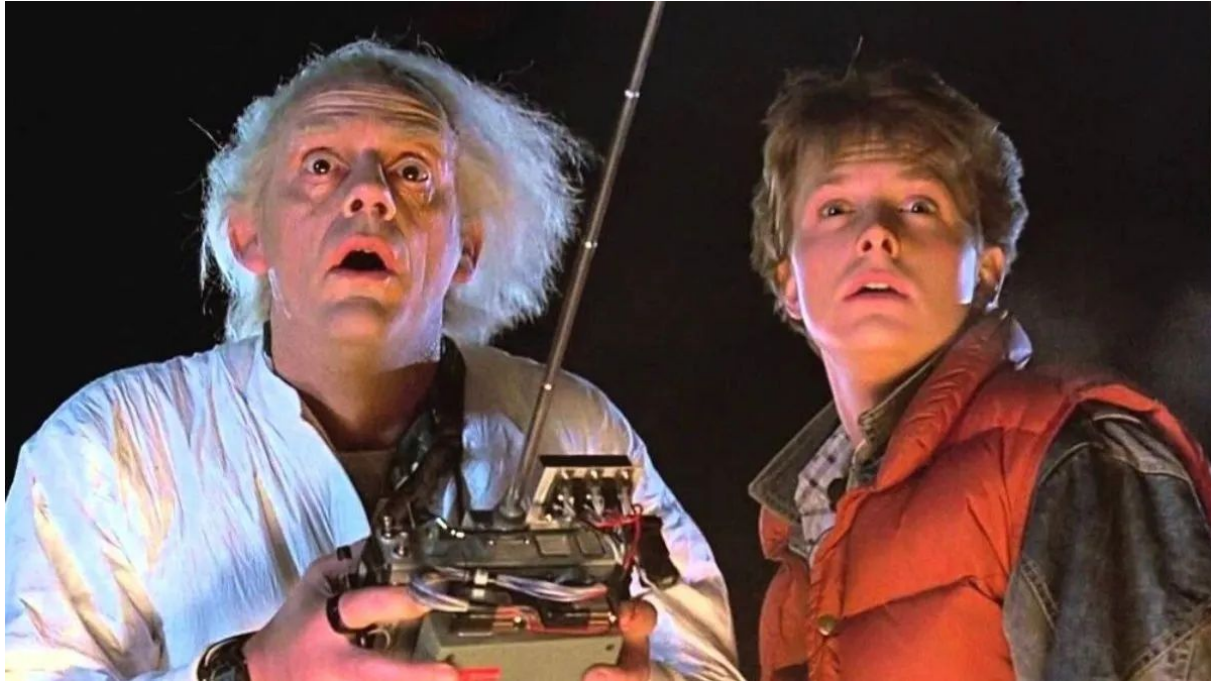
$$= (1-\gamma)^{-1} \mathbb{E}_{\pi} [\mathbb{E}_{\gamma} [R_{\tau+t}] | S_t = s]$$

$$E_{\gamma}[R_{\tau+t}] = E_{\gamma}[R_{\tau+t} 1_{\{\tau \leq n\}}] + P(\tau > n) E_{\gamma}[R_{\tau+t} | \tau > n]$$

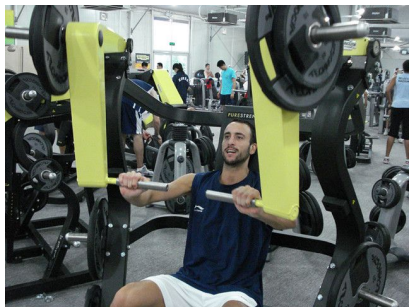
$$G_t^{(n)} = E_{\pi}[E_{\gamma}[R_{\tau+t}] | \mathcal{F}_n]$$

$$= E_{\gamma}[R_{\tau+t} 1_{\{\tau \leq n\}}] + \gamma^n v(S_{t+n})$$

PERO TD(LAMBDA) REQUIERE DE TODO EL EPISODIO HASTA EL FINAL...TENGO QUE MIRAR AL FUTURO...



ELIGIBILITY TRACES



¿qué explica mayormente el campeonato?

- El entrenamiento (frecuencia)
- Las vacaciones (cuán recientemente visité ese estado)

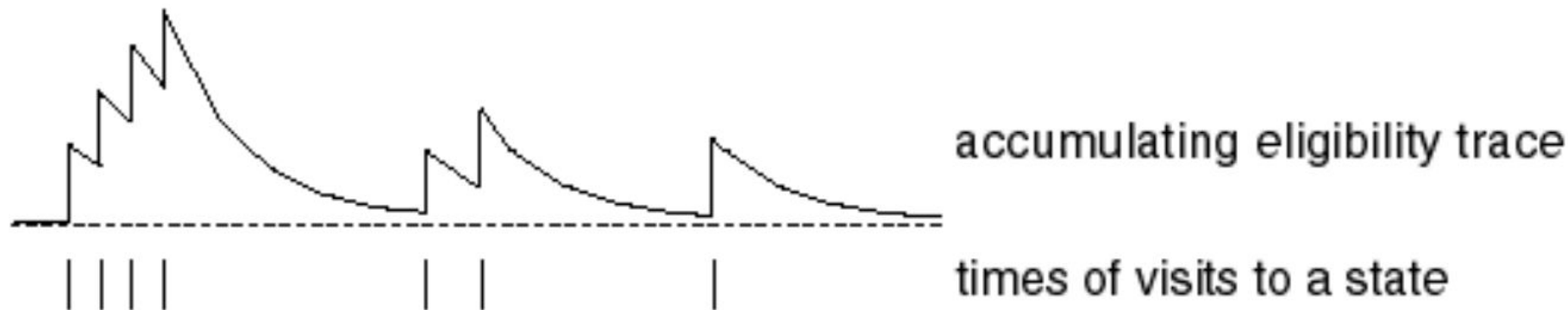
¿cómo combinar ambas?



ELIGIBILITY TRACES

$$E_0(s) = 0$$

$$E_t(s) = \gamma\lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$$



$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$
$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

FOWARD-VIEW \sim BACKWARD-VIEW

Theorem

The sum of offline updates is identical for forward-view and backward-view TD(λ)

$$\sum_{t=1}^T \alpha \delta_t E_t(s) = \sum_{t=1}^T \alpha \left(G_t^\lambda - V(S_t) \right) \mathbf{1}(S_t = s)$$

EL LADO OSCURO...



Off-policy Learning with Eligibility Traces: A Survey

<https://hal.inria.fr/hal-00644516/PDF/jmlr.pdf>

Convergence Results for Single-Step On-Policy Reinforcement-Learning Algorithms

<https://link.springer.com/content/pdf/10.1023/A:1007678930559.pdf>

