

CARTA

doi:10.1038/natureza14236

Control a nivel humano a través del aprendizaje de refuerzo profundo

Volodymyr Mnih¹*, Koray Kavukcuoglu¹*, David Plata¹*, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martín Riedmiller¹, Andreas K. Fidjeland¹, Jorge Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen Rey¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ y Demis Hassabis¹

La teoría del aprendizaje por refuerzo proporciona una explicación normativa¹, profundamente arraigada en la psicología y neurocientífica² perspectivas sobre el comportamiento animal, de cómo los agentes pueden optimizar su control de un entorno. Estos para generalizar la experiencia pasada a situaciones nuevas. Sorprendentemente, los humanos y otros animales parecen resolver este problema a través de una combinación armoniosa de aprendizaje por refuerzo y sistemas jerárquicos de procesamiento sensorial^{4,5}, el primero evidenciado por una gran cantidad de datos neuronales que revelan paralelos notables entre las señales fisiológicas emitidas por las neuronas dopaminérgicas y los algoritmos de aprendizaje de refuerzo de diferencia temporal³. Si bien los agentes de aprendizaje por refuerzo han logrado algunos éxitos en una variedad de dominios⁶⁻⁸, su aplicabilidad se ha limitado previamente a dominios en los que las características útiles se pueden crear a mano, o a dominios con espacios de estado de baja dimensión completamente observados. Aquí usamos avances recientes en el entrenamiento de redes neuronales profundas⁹⁻¹¹ para desarrollar un agente artificial novedoso, denominado red Q profunda, que puede aprender políticas exitosas directamente de entradas sensoriales de alta dimensión utilizando aprendizaje de refuerzo de extremo a extremo. Probamos este agente en el desafiante dominio de los juegos clásicos de Atari 2600¹². Demostramos que el agente de la red Q profunda, al recibir solo los píxeles y la puntuación del juego como entradas, pudo superar el rendimiento de todos los algoritmos anteriores y alcanzar un nivel comparable al de un probador profesional de juegos humanos en un conjunto de 49 juegos, usando el mismo algoritmo, arquitectura de red e hiperparámetros. Este trabajo cierra la brecha entre las acciones y las entradas sensoriales de alta dimensión, lo que da como resultado el primer agente artificial que es capaz de aprender a sobresalir en una amplia gama de tareas desafiantes.

Nos propusimos crear un algoritmo único que pudiera desarrollar una amplia gama de competencias en una variedad de tareas desafiantes, un objetivo central de la inteligencia artificial general¹³ que ha eludido esfuerzos anteriores^{8,14,15}. Para lograr esto, desarrollamos un nuevo agente, una red profunda Q (DQN), que es capaz de combinar el aprendizaje por refuerzo con una clase de red neuronal artificial¹⁶, conocidas como redes neuronales profundas. En particular, los avances recientes en redes neuronales profundas⁹⁻¹¹, en el que se utilizan varias capas de nodos para construir representaciones progresivamente más abstractas de los datos, han hecho posible que las redes neuronales artificiales aprendan conceptos como categorías de objetos directamente a partir de datos sensoriales sin procesar. Usamos una arquitectura particularmente exitosa, la red convolucional profunda¹⁷, que utiliza capas jerárquicas de filtros convolucionales en mosaico para imitar los efectos de los campos receptivos, inspirado en el trabajo seminal de Hubel y Wiesel sobre el procesamiento de avance en la corteza visual temprana¹⁸—explotando así las correlaciones espaciales locales presentes en las imágenes, y construyendo solidez a las transformaciones naturales como cambios de punto de vista o escala.

Consideramos tareas en las que el agente interactúa con un entorno a través de una secuencia de observaciones, acciones y recompensas. El objetivo del

agente es seleccionar acciones de una manera que maximice la recompensa futura acumulada. Más formalmente, usamos una red neuronal convolucional profunda para aproximar la función de valor de acción óptima

$$q(s,a) \approx \max_{a'} r + \gamma \max_{s'} V(s')$$

que es la suma máxima de recompensas descontada por γ cada paso de tiempo, alcanzable por una política de comportamiento $\pi(a|s)$, después de hacer una observación (s) y tomando una acción (a) (ver Métodos)¹⁹.

Se sabe que el aprendizaje por refuerzo es inestable o incluso diverge cuando se usa un aproximador de función no lineal, como una red neuronal, para representar el valor de la acción (también conocido como Q -función)²⁰. Esta inestabilidad tiene varias causas: las correlaciones presentes en la secuencia de observaciones, el hecho de que pequeñas actualizaciones de Q pueden cambiar significativamente la política y, por lo tanto, cambiar la distribución de datos y las correlaciones entre los valores de acción (q) y los valores objetivos²¹.

Abordamos estas inestabilidades con una variante novedosa de Q -learning, que utiliza dos ideas clave. Primero, usamos un mecanismo de inspiración biológica llamado repetición de experiencia²¹⁻²³ que aleatoriza los datos, eliminando así las correlaciones en la secuencia de observaciones y suavizando los cambios en la distribución de datos (consulte los detalles a continuación). En segundo lugar, utilizamos una actualización iterativa que ajusta los valores de acción (q) hacia valores objetivo que solo se actualizan periódicamente, reduciendo así las correlaciones con el objetivo.

Si bien existen otros métodos estables para entrenar redes neuronales en el entorno de aprendizaje por refuerzo, como la iteración Q ajustada neuronal²⁴, estos métodos implican el entrenamiento repetido de redes durante cientos de iteraciones. En consecuencia, estos métodos, a diferencia de nuestro algoritmo, son demasiado ineficientes para ser utilizados con éxito con grandes redes neuronales. Parametrizamos una función de valor aproximado $Q(s,a;h)$ utilizando la red neuronal convolucional profunda que se muestra en la Fig. 1, en la que se describen los parámetros (es decir, los pesos) de la red Q en la iteración i . Para realizar la repetición de la experiencia, almacenamos las experiencias del agente (s_i, a_i, r_i, s_{i+1}) en cada paso de tiempo en un conjunto de datos $D = \{(s_1, a_1, r_1, s_2), \dots, (s_{T-1}, a_{T-1}, r_{T-1}, s_T)\}$. Durante el aprendizaje, aplicamos actualizaciones de Q -learning, sobre muestras (o mini-batches) de experiencia (s, a, r, s') , extraídas uniformemente al azar del grupo de muestras almacenadas. La actualización de Q -learning en la iteración i utiliza la siguiente función de pérdida:

$$L(Q, \pi, D) = \mathbb{E}_{(s, a, r, s') \sim D} \left[\max_{a'} Q(s, a; h_i) - \mathbb{E}_{a' \sim \pi(\cdot | s)} [Q(s, a'; h_i)] \right]$$

en el cual \mathbb{E} es el factor de descuento que determina el horizonte del agente, h son los parámetros de la red Q en la iteración i , π es la red Q utilizada para calcular el objetivo en la iteración i . Los parámetros de trabajo h solo se actualizan con los parámetros de la red Q (h_i) cada T pasos y se mantienen fijos entre actualizaciones individuales (ver Métodos).

Para evaluar nuestro agente DQN, aprovechamos la plataforma Atari 2600, que ofrece una diversa gama de tareas (norte 549) diseñado para ser

¹Google DeepMind, 5 New Street Square, Londres EC4A 3TW, Reino Unido.

* Estos autores contribuyeron igualmente a este trabajo.

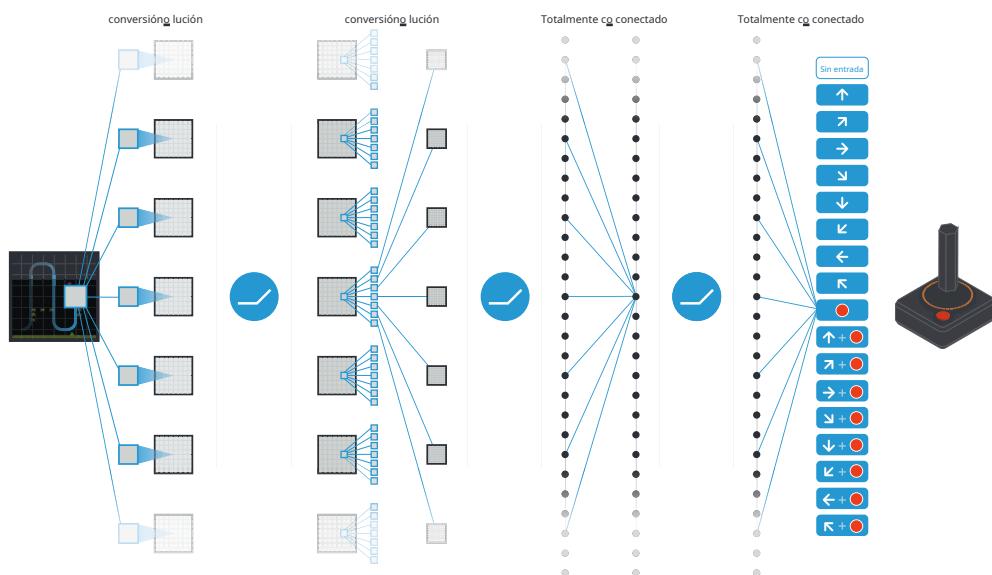


Figura 1 | Ilustración esquemática de la red neuronal convolucional. Los detalles de la arquitectura se explican en los Métodos. La entrada a la red neuronal consiste en un 8438434 imagen producida por el mapa de preprocesamiento, seguido de tres capas convolucionales (nota: línea azul serpenteante

difícil y atractivo para los jugadores humanos. Usamos la misma arquitectura de red, valores de hiperparámetros (consulte la Tabla 1 de datos ampliados) y un procedimiento de aprendizaje en todo momento (tomando datos de alta dimensión (210|160 video en color a 60 Hz) como entrada) para demostrar que nuestro enfoque aprende de manera sólida políticas exitosas en una variedad de juegos basados únicamente en entradas sensoriales con solo un conocimiento previo muy mínimo (es decir, solo los datos de entrada eran imágenes visuales y la cantidad de acciones disponibles en cada juego, pero no sus correspondencias; ver Métodos). En particular, nuestro método fue capaz de entrenar grandes redes neuronales utilizando una señal de aprendizaje de refuerzo y un descenso de gradiente estocástico de una manera estable ilustrada por la evolución temporal de dos índices de aprendizaje (la puntuación promedio por episodio del agente y los valores Q promedio pronosticados; consulte la Fig.

symboliza el deslizamiento de cada filtro a través de la imagen de entrada) y dos capas completamente conectadas con una sola salida para cada acción válida. A cada capa oculta le sigue una no linealidad del rectificador (es decir, $\max(0, x)$).

Comparamos DQN con los métodos de mejor desempeño de la literatura de aprendizaje por refuerzo en los 49 juegos donde los resultados estaban disponibles.^{12,15} Además de los agentes aprendidos, también informamos puntuajes para un probador profesional de juegos humanos que juega bajo condiciones controladas y una política que selecciona acciones uniformemente al azar (Tabla de datos extendidos 2 y Fig. 3, indicada por 100% (humano) y 0% (aleatorio) en el eje; ver Métodos). Nuestro método DQN supera a los mejores métodos de aprendizaje por refuerzo existentes en 43 de los juegos sin incorporar ninguno de los conocimientos previos adicionales sobre los juegos de Atari 2600 utilizados por otros enfoques (por ejemplo, refs 12, 15). Además, nuestro agente de DQN se desempeñó a un nivel comparable al de un probador profesional de juegos humanos en el conjunto de 49 juegos, logrando más del 75 % de la puntuación humana en más de la mitad de los juegos (29 juegos);

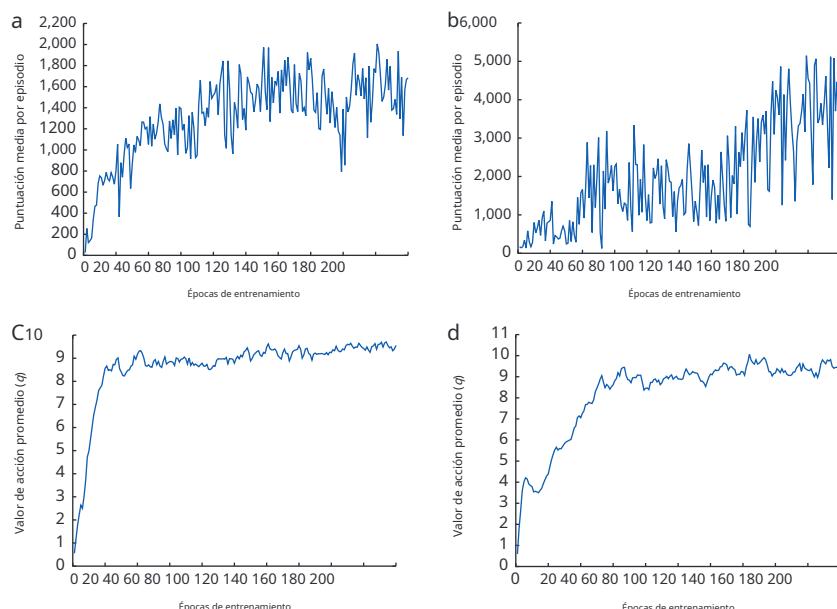


Figura 2 | Curvas de entrenamiento que rastrean el puntaje promedio del agente y el valor de acción pronosticado promedio. a, Cada punto es la puntuación promedio lograda por episodio después de que el agente se ejecuta con una política codiciosa ($\epsilon=0.05$) para cuadros de 520k en Space Invaders. b, Puntaje promedio logrado por episodio para Seaquest. c, Valor de acción pronosticado promedio en un conjunto de estados retenidos en Space Invaders. Cada punto

en la curva está el promedio del valor de acción Q calculado sobre el conjunto de estados retenidos. Tenga en cuenta que los valores Q se escalan debido al recorte de recompensas (ver Métodos). d, Valor de acción promedio previsto en Seaquest. Ver Discusión Complementaria para más detalles.

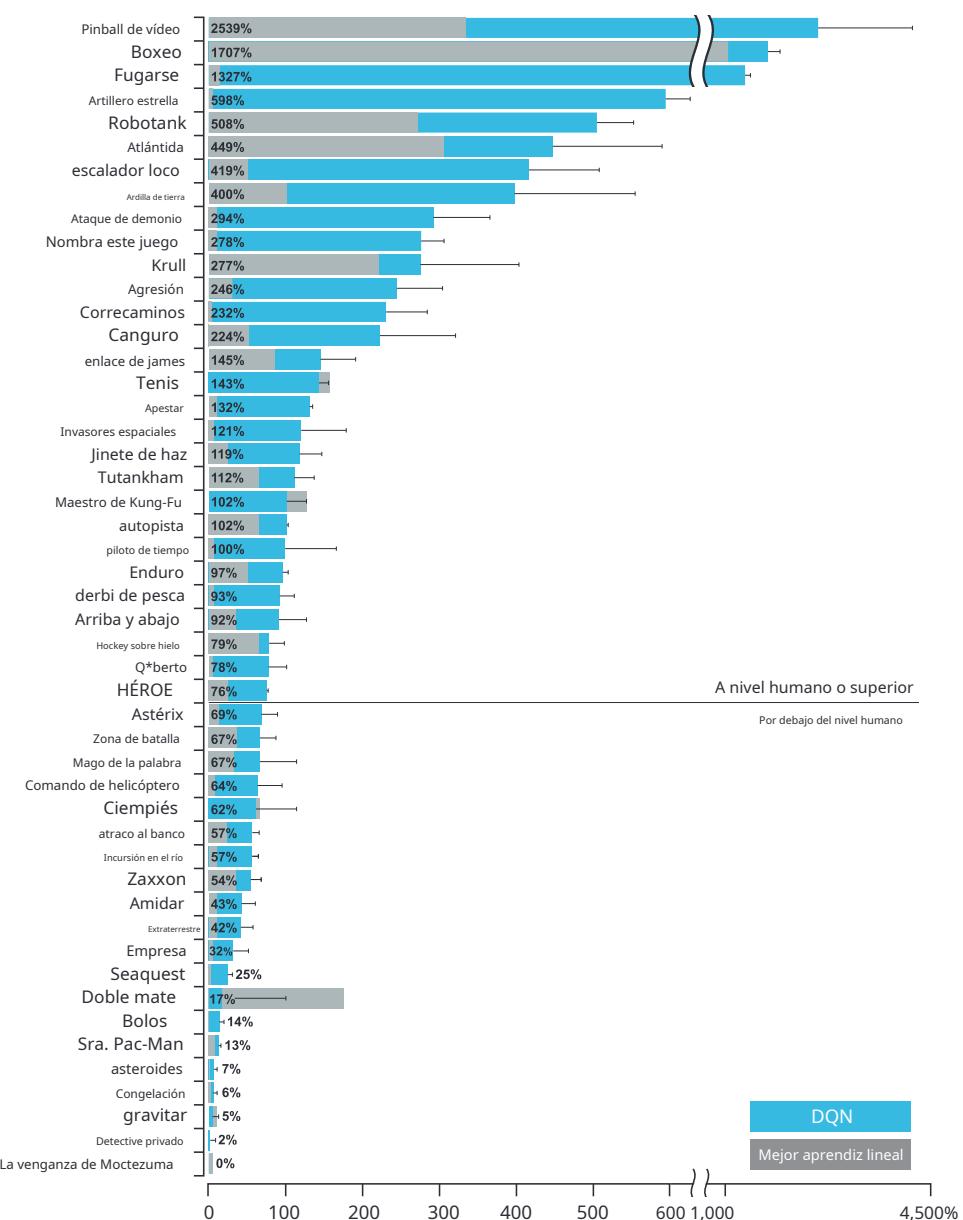


Figura 3 | Comparación del agente DQN con los mejores métodos de aprendizaje por refuerzo en la literatura. El rendimiento de DQN está normalizado con respecto a un probador de juegos humanos profesionales (es decir, nivel 100 %) y juego aleatorio (es decir, nivel 0 %). Tenga en cuenta que el rendimiento normalizado de DQN, expresado como porcentaje, se calcula como: 1003(Puntuación DQN2puntuación de juego aleatorio)/(puntuación humana2puntuación de juego aleatorio). Se puede ver que DQN

superia a los métodos de la competencia (consulte también la Tabla 2 de datos ampliados) en casi todos los juegos, y se desempeña a un nivel que es ampliamente comparable o superior al de un evaluador de juegos humano profesional (es decir, operacionalizado como un nivel del 75 % o superior) en la mayoría de los juegos. La salida de audio se deshabilitó tanto para jugadores humanos como para agentes. Las barras de error indican sd en los 30 episodios de evaluación, comenzando con diferentes condiciones iniciales.

consulte la Fig. 3, Discusión complementaria y Tabla de datos ampliados 2). En simulaciones adicionales (consulte la discusión complementaria y las tablas de datos extendidos 3 y 4), demostramos la importancia de los componentes centrales individuales del DQNagent (la memoria de reproducción, la red Q de destino separada y la arquitectura de red convolucional profunda) al desactivarlos y demostrar los efectos perjudiciales en el rendimiento.

A continuación, examinamos las representaciones aprendidas por DQN que respaldaron el desempeño exitoso del agente en el contexto del juego Space Invaders (consulte el Video complementario 1 para ver una demostración del desempeño de DQN), mediante el uso de una técnica desarrollada para la visualización de datos de alta dimensión llamada 't-SNE'²⁵(Figura 4). Como era de esperar, el algoritmo t-SNE tiende a mapear la representación DQN de estados perceptualmente similares en puntos cercanos. Curiosamente, también encontramos instancias en las que el algoritmo t-SNE generó incrustaciones similares para representaciones DQN de estados que están cerca en términos de recompensa esperada pero

perceptualmente diferentes (Fig. 4, abajo a la derecha, arriba a la izquierda y en el medio), de acuerdo con la noción de que la red puede aprender representaciones que respaldan el comportamiento adaptativo a partir de entradas sensoriales de alta dimensión. Además, también mostramos que las representaciones aprendidas por DQN pueden generalizarse a datos generados a partir de políticas distintas a las suyas, en simulaciones donde presentamos como entrada a los estados del juego de red experimentados durante el juego humano y de agentes, registramos las representaciones de la última capa oculta., y visualizó las incrustaciones generadas por el algoritmo t-SNE (Datos extendidos, Fig. 1 y Discusión complementaria). ExtendedData Fig. 2 proporciona una ilustración adicional de cómo las representaciones aprendidas por DQN le permiten predecir con precisión los valores de estado y acción.

Vale la pena señalar que los juegos en los que DQN sobresale son muy variados en su naturaleza, desde juegos de disparos de desplazamiento lateral (River Raid) hasta juegos de boxeo (Boxing) y juegos de carreras de autos tridimensionales (Enduro).

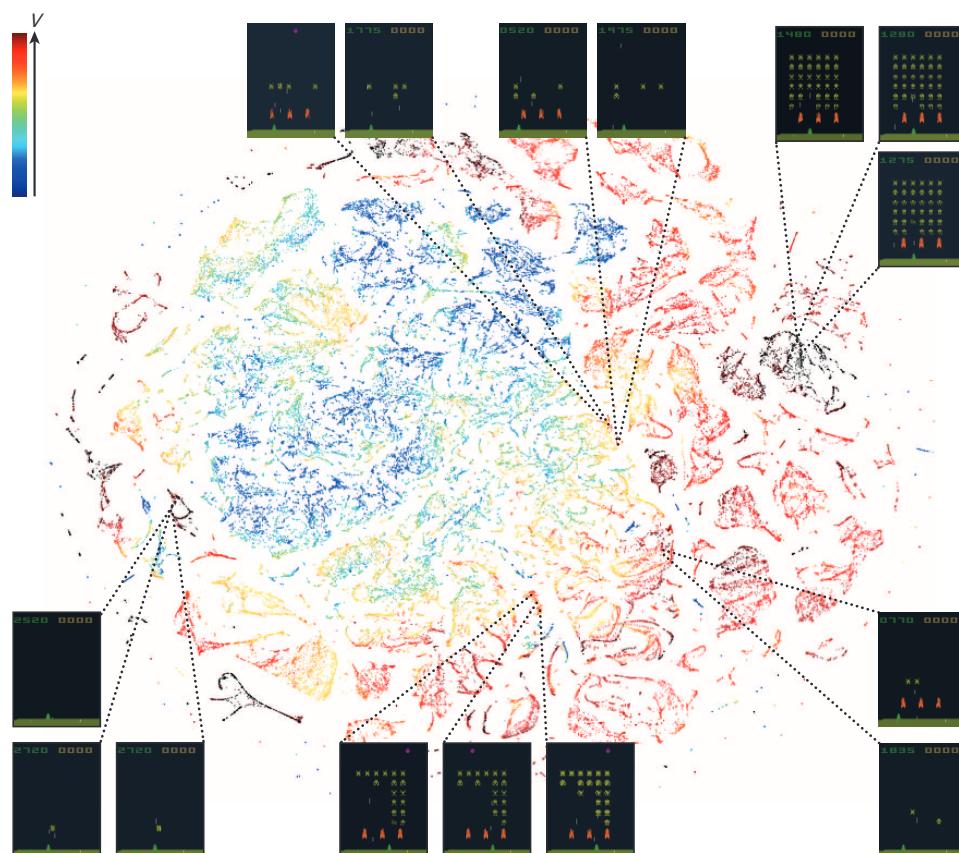


Figura 4 | Incrustación t-SNE bidimensional de las representaciones en la última capa oculta asignada por DQN a los estados del juego experimentados al jugar a Space Invaders. La trama se generó dejando que el agente de DQN jugara durante 2 h de tiempo real de juego y ejecutando el algoritmo t-SNE...en las últimas representaciones de capa oculta asignadas por DQN a cada estado de juego experimentado. Los puntos se colorean de acuerdo con los valores de estado (v , recompensa máxima esperada de un estado) predicha por DQN para los estados de juego correspondientes (que van desde el rojo oscuro (más alto) a azul oscuro (más bajo)). Se muestran las capturas de pantalla correspondientes a un número seleccionado de puntos. El agente DQN

predice valores de estado altos tanto para pantallas completas (capturas de pantalla en la parte superior derecha) como casi completas (capturas de pantalla en la parte inferior izquierda) porque ha aprendido que completar una pantalla conduce a una nueva pantalla llena de naves enemigas. A las pantallas parcialmente completadas (capturas de pantalla inferiores) se les asignan valores de estado más bajos porque hay una recompensa menos inmediata disponible. Las pantallas que se muestran abajo a la derecha, arriba a la izquierda y en el medio son menos similares a la percepción que los otros ejemplos, pero todavía están asignadas a representaciones cercanas y valores similares porque los bunkers naranjas no tienen gran importancia cerca del final de un nivel. Con permiso de Square Enix Limited.

De hecho, en ciertos juegos, DQN puede descubrir una estrategia relativamente a largo plazo (por ejemplo, Breakout: el agente aprende la estrategia óptima, que consiste en cavar primero un túnel alrededor del costado de la pared que permite que la pelota se envíe alrededor del vóter a destruir una gran cantidad de bloques; vea el video complementario 2 para ver una ilustración del desarrollo del desempeño de DQN durante el transcurso del entrenamiento).

En este trabajo, demostramos que una sola arquitectura puede aprender con éxito políticas de control en una variedad de entornos diferentes con un conocimiento previo mínimo, recibiendo solo los píxeles y la puntuación del juego como entradas, y usando el mismo algoritmo, arquitectura de red e hiperparámetros en cada juego, con acceso solo a las entradas que tendría un jugador humano. A diferencia del trabajo anterior^{24,26}, nuestro enfoque incorpora el aprendizaje por refuerzo de "extremo a extremo" que utiliza la recompensa para dar forma continua a las representaciones dentro de la red convolucional hacia características sobresalientes del entorno que facilitan la estimación del valor. Este principio se basa en la evidencia neurobiológica de que las señales de recompensa durante el aprendizaje perceptivo pueden influir en las características de las representaciones dentro de la corteza visual de los primates.^{27,28} En particular, la integración exitosa del aprendizaje por refuerzo con arquitecturas de redes profundas dependía de manera crítica de nuestra incorporación de un algoritmo de reproducción.²¹⁻²³ que implica el almacenamiento y la representación de transiciones experimentadas recientemente. Evidencia convergente sugiere que el hipocampo puede soportar el

realización de tal proceso en el cerebro de los mamíferos, con la reactivación comprimida en el tiempo de trayectorias experimentadas recientemente durante períodos fuera de línea^{21,22} (por ejemplo, descanso despierto) proporcionando un mecanismo putativo por el cual las funciones de valor pueden actualizarse de manera eficiente a través de interacciones con los ganglios basales²². En el futuro, será importante explorar el uso potencial de sesgar el contenido de la repetición de la experiencia hacia eventos destacados, un fenómeno que caracteriza la repetición del hipocampo observada empíricamente.²⁹ y se relaciona con la noción de 'barrido priorizado'³⁰ en el aprendizaje por refuerzo. En conjunto, nuestro trabajo ilustra el poder de aprovechar las técnicas de aprendizaje automático de última generación con mecanismos inspirados biológicamente para crear agentes que sean capaces de aprender a dominar una amplia gama de tareas desafiantes.

Contenido en línea Los métodos, junto con cualquier elemento adicional de visualización de datos extendidos y datos de origen, están disponibles en la versión en línea del documento; las referencias únicas a estas secciones aparecen solo en el documento en línea.

Recibido el 10 de julio de 2014; aceptado el 16 de enero de 2015.

1. Sutton, R. y Barto, A. Aprendizaje por refuerzo: una introducción (Presa del MIT, 1998).
2. Thorndike, E. L. Inteligencia animal: estudios experimentales (Macmillan, 1911).
3. Schultz, W., Dayan, P. & Montague, P. R. Un sustrato neuronal de predicción y recompensa. *Ciencias* 275, 1593-1599 (1997).
4. Serre, T., Wolf, L. & Poggio, T. Reconocimiento de objetos con características inspiradas en la corteza visual. *Proc. IEEE computar Soc. Conf. computar Vis. Patrón. reconocer* 994-1000 (2005).
5. Fukushima, K. Neocognitron: un modelo de red neuronal autoorganizada para un mecanismo de reconocimiento de patrones que no se ve afectado por el cambio de posición. *Biol. Cybern.* 36, 193-202 (1980).

6. Tesauro, G. Aprendizaje de diferencias temporales y TD-Gammon.común MCA38, 58–68 (1995).
7. Riedmiller, M., Gabel, T., Hafner, R. y Lange, S. Aprendizaje por refuerzo para fútbol robótico. *Auton. robots*27, 55–73 (2009).
8. Diuk, C., Cohen, A. & Littman, ML Una representación orientada a objetos para un aprendizaje reforzado eficiente.proc. En t. Conf.Mach. Aprender.240–247 (2008).
9. Bengio, Y. Aprendizaje de arquitecturas profundas para IA.Fundamentos y Tendencias en Machine Learning2,1-127 (2009).
10. Krizhevsky, A., Sutskever, I. & Hinton, G. Clasificación de ImageNet con redes neuronales convolucionales profundas. *Adv.Neural Inf.Process.Syst.*25, 1106–1114 (2012).
11. Hinton, GE & Salakhutdinov, RR Reducción de la dimensionalidad de los datos con redes neuronales. *Ciencias*313, 504–507 (2006).
12. Bellemare, MG, Naddaf, Y., Veness, J. & Bowling, M. El entorno de aprendizaje arcade: una plataforma de evaluación para agentes generales. *J. Artif. Intel. Res.*47, 253–279 (2013).
13. Legg, S. & Hutter, M. Inteligencia universal: una definición de inteligencia artificial. *Minds Mach.*17, 391–444 (2007).
14. Genesereth, M., Love, N. & Pell, B. Juego general: descripción general de la competencia AAAI.IA Mag.26, 62–72 (2005).
15. Bellemare, MG, Veness, J. & Bowling, M. Investigación de la conciencia de contingencia utilizando juegos Atari 2600.proc. Conf. AAAI. Artefacto Intel.864–871 (2012).
16. McClelland, JL, Rumelhart, DE y Grupo, TPRProcesamiento distribuido en paralelo: exploraciones en la microestructura de la cognición (Prensa del MIT, 1986).
17. LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. Aprendizaje basado en gradiente aplicado al reconocimiento de documentos.proc. IEEE86, 2278–2324 (1998).
18. Hubel, DH y Wiesel, TN Forma y disposición de las columnas en la corteza estriada del gato. *J. Physiol.*165, 559–568 (1963).
19. Watkins, CJ y Dayan, P. Q-learning. *Mach. Aprender.*8, 279–292 (1992).
20. Tsitsiklis, J. & Roy, BV Un análisis del aprendizaje de diferencias temporales con aproximación de funciones. *IEEE Trans. Aparato mecánico. Contr.*42, 674–690 (1997).
21. McClelland, JL, McNaughton, BL y O'Reilly, RC Por qué existen sistemas de aprendizaje complementarios en el hipocampo y el neocórtex: percepciones de los éxitos y fracasos de los modelos conexiónistas de aprendizaje y memoria. *psicol. Rvdo.*102, 419–457 (1995).
22. O'Neill, J., Pleydell-Bouverie, B., Dupret, D. & Csicsvari, J. Tócala de nuevo: reactivación de la experiencia de vigilia y la memoria. *Tendencias Neurosci.*33, 220–229 (2010).
23. Lin, L.-J. Aprendizaje por refuerzo para robots mediante redes neuronales. Informe Técnico, Documento DTIC (1993).
24. Riedmiller, M. Neural equipado Q iteración: primeras experiencias con un método de aprendizaje de refuerzo neuronal eficiente en datos. *Mach. Aprender.: ECML*,3720,317–328 (Springer, 2005).
25. Van der Maaten, LJP & Hinton, GE Visualización de datos de alta dimensión usando t-SNE. *J. Mach. Aprender. Res.*9, 2579–2605 (2008).
26. Lange, S. & Riedmiller, M. Redes neuronales de autocodificador profundo en el aprendizaje por refuerzo.proc. En t. *Jt. Conf. Neural. Neto.*1–8 (2010).
27. Ley, C.-T. & Gold, JI El aprendizaje por refuerzo puede dar cuenta del aprendizaje asociativo y perceptivo en una tarea de decisión visual. *Naturaleza Neurosci.*12, 655 (2009).
28. Sigala, N. & Logothetis, NK Las formas de categorización visual presentan selectividad en la corteza temporal de los primates. *Naturaleza*415, 318–320 (2002).
29. Bendor, D. & Wilson, MA Sesgo del contenido de la reproducción del hipocampo durante el sueño. *Naturaleza Neurosci.*15, 1439–1444 (2012).
30. Moore, A. & Atkeson, C. Barrido priorizado: aprendizaje por refuerzo con menos datos y menos tiempo real. *Mach. Aprender.*13, 103–130 (1993).

Información suplementariaestá disponible en la versión en línea del documento.

AgradecimientosAgradecemos a G. Hinton, P. Dayan y M. Bowling por las discusiones, a A. Cain y J. Keene por trabajar en las imágenes, a K. Keller y P. Rogers por su ayuda con las imágenes, a G. Wayne por sus comentarios sobre una versión anterior. versión del manuscrito, y al resto del equipo de DeepMind por su apoyo, ideas y aliento.

Contribuciones de autorVM, KK, DS, JV, MGB, MR, AG, DW, SL y DH conceptualizaron el problema y el marco técnico. VM, KK, AAR y DS desarrollaron y probaron los algoritmos. JV, SP, CB, AAR, MGB, IA, AKF, GO y AS crearon la plataforma de pruebas. KK, HK, SL y D.H. manejó el proyecto. KK, DK, DH, VM, DS, AG, AAR, JV y MGB escribieron el artículo.

Información del autorLa información sobre reimpresiones y permisos está disponible en www.nature.com/reprints. Los autores declaran no tener intereses financieros en competencia. Los lectores pueden comentar sobre la versión en línea del documento. La correspondencia y las solicitudes de materiales deben dirigirse a KK (korayk@google.com) o DH (demishassabis@google.com).

MÉTODOS

Preprocesamiento.Trabajar directamente con marcos Raw Atari 2600, que son 2103Las imágenes de 160 píxeles con una paleta de 128 colores pueden ser exigentes en términos de requisitos de cómputo y memoria. Aplicamos un paso básico de preprocesamiento destinado a reducir la dimensionalidad de entrada y lidiar con algunos artefactos del emulador Atari 2600. Primero, para codificar un solo cuadro, tomamos el valor máximo para cada valor de color de pixel sobre el cuadro que se está codificando y el cuadro anterior. Esto fue necesario para eliminar el parpadeo que está presente en los juegos donde algunos objetos aparecen solo en cuadros pares mientras que otros objetos aparecen solo en cuadros impares, un artefacto causado por la cantidad limitada de sprites que Atari 2600 puede mostrar a la vez. En segundo lugar, extraemos el canal Y, también conocido como luminancia, del marco RGB y lo cambiamos de escala a 84384. La función wfromalgorith1 descrito a continuación aplica este preprocesamiento almetrototogramas más recientes y los apila para producir la entrada a la función Q, en la quemetro54, aunque el algoritmo es robusto a diferentes valores de m (por ejemplo, 3 o 5).

Disponibilidad de código.Se puede acceder al código fuente en <https://sites.google.com/a/deepmind.com/dqn> solo para usos no comerciales.

Modelo de arquitectura.Hay varias formas posibles de parametrizar Q utilizando una red neuronal. Debido a que Q asigna pares de historia-acción a estimaciones escalares de su valor Q, algunos enfoques anteriores han utilizado la historia y la acción como entradas a la red neuronal.^{24,26} El principal inconveniente de este tipo de arquitectura es que se requiere un paso adelante separado para calcular el valor Q de cada acción, lo que da como resultado un costo que escala linealmente con el número de acciones. En su lugar, usamos una arquitectura en la que hay una unidad de salida separada para cada acción posible, y solo la representación del estado es una entrada para la red neuronal. Las salidas corresponden a los valores Q predichos de las acciones individuales para el estado de entrada. La principal ventaja de este tipo de arquitectura es la capacidad de calcular los valores Q para todas las acciones posibles en un estado dado con un solo paso directo a través de la red.

La arquitectura exacta, que se muestra esquemáticamente en la Fig. 1, es la siguiente. La entrada a la red neuronal consiste en un 8438434 imagen producida por el mapa de preprocesamientow.La primera capa oculta convoluciona 32 filtros de 838 con zancada 4 con la imagen de entrada y aplica una no linealidad rectificadora^{31,32}. La segunda capa oculta convoluciona 64 filtros de 434 con paso 2, seguido de nuevo por una no linealidad rectificadora. A esto le sigue una tercera capa convolucional que convoluciona 64 filtros de 333 con paso 1 seguido de un rectificador. La capa oculta final está totalmente conectada y consta de 512 unidades rectificadoras. La capa de salida es una capa lineal completamente conectada con una sola salida para cada acción válida. El número de acciones válidas varió entre 4 y 18 en los juegos considerados.

Detalles de entrenamiento.Realizamos experimentos en 49 juegos Atari 2600 donde los resultados estaban disponibles para todos los demás métodos comparables^{12,15}. Se entrenó una red diferente en cada juego: se usó la misma arquitectura de red, algoritmo de aprendizaje y configuración de hiperparámetros (consulte la Tabla de datos ampliados 1) en todos los juegos, lo que demuestra que nuestro enfoque es lo suficientemente sólido para funcionar en una variedad de juegos mientras incorpora solo un conocimiento previo mínimo (consulte a continuación). Mientras evaluamos a nuestros agentes en juegos no modificados , hicimos un cambio en la estructura de recompensas de los juegos solo durante el entrenamiento. Como la escala de puntajes varía mucho de un juego a otro, recortamos todas las recompensas positivas en 1 y todas las recompensas negativas en 21, dejando 0 recompensas sin cambios. Recortar las recompensas de esta manera limita la escala de los errores derivados y facilita el uso de la misma tasa de aprendizaje en varios juegos. Al mismo tiempo, podría afectar al rendimiento de nuestro agente ya que no puede diferenciar entre recompensas de distinta magnitud. Para los juegos donde hay un contador de vidas, el emulador Atari 2600 también envía la cantidad de vidas que quedan en el juego, que luego se usa para marcar el final de un episodio durante el entrenamiento.

En estos experimentos, usamos el RMSProp (ver http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf) algoritmo con minibatches de tamaño 32. La política de comportamiento durante el entrenamiento fuem-iocidioso comnirecodido linealmente de 1,0 a 0,1 durante el primer millón de fotogramas y fijado en 0,1 a partir de entonces. Entrenamos para un total de 50 millones de fotogramas (es decir, alrededor de 38 días de experiencia de juego en total) y utilizamos una memoria de repetición de 1 millón de fotogramas más recientes.

Siguiendo los enfoques anteriores para jugar juegos Atari 2600, también usamos una técnica simple de salto de marco¹⁵. Más precisamente, el agente ve y selecciona acciones en cadakel fotograma en lugar de cada fotograma, y su última acción se repite en los fotogramas omitidos. Debido a que ejecutar el emulador un paso hacia adelante requiere mucho menos cálculo que hacer que el agente seleccione una acción, esta técnica permite que el agente juegue aproximadamentekveces más juegos sin aumentar significativamente el tiempo de ejecución. Usamosk54 para todos los juegos.

Los valores de todos los hiperparámetros y parámetros de optimización se seleccionaron realizando una búsqueda informal en los juegos Pong, Breakout, Seaquest, Space Invaders y Beam Rider. No realizamos una búsqueda sistemática en cuadrícula debido al alto costo computacional. Luego, estos parámetros se mantuvieron fijos en todos los demás juegos. Los valores y las descripciones de todos los hiperparámetros se proporcionan en la Tabla 1 de datos ampliados.

Nuestra configuración experimental consiste en utilizar el siguiente conocimiento previo mínimo: que los datos de entrada consistían en imágenes visuales (lo que motiva nuestro uso de una red profunda convolucional), la puntuación específica del juego (sin modificaciones), el número de acciones, aunque no sus correspondencias (por ejemplo, especificación del 'botón arriba') y el recuento de vida.

Procedimiento de evaluación.Los agentes capacitados se evaluaron jugando cada juego 30 veces durante un máximo de 5 minutos cada vez con diferentes condiciones aleatorias iniciales ("noop"; consulte la Tabla 1 de datos ampliados) y umni-política codicosa commi50.05. Este procedimiento se adopta para minimizar la posibilidad de sobreajuste durante la evaluación. El agente aleatorio sirvió como comparación de referencia y eligió una acción aleatoria a 10 Hz, que es cada sexto cuadro, repitiendo su última acción en los cuadros intermedios. 10 Hz es aproximadamente lo más rápido que un jugador humano puede seleccionar el botón 'disparar', y configurar el agente aleatorio en esta frecuencia evita puntajes de referencia falsos en un puñado de juegos. También evaluamos el desempeño de un agente aleatorio que seleccionó una acción a 60 Hz (es decir, cada cuadro). Esto tuvo un efecto mínimo: cambió el rendimiento normalizado de DQN en más del 5 % en solo seis juegos (Boxing, Breakout, Crazy Climber, Demon Attack, Krull y Robotank), y en todos estos juegos DQN superó al experto humano por un margen considerable. .

El evaluador humano profesional utilizó el mismo motor de emulación que los agentes y jugó en condiciones controladas. El probador humano no podía pausar, guardar o recargar juegos. Como en el entorno Atari 2600 original, el emulador se ejecutaba a 60 Hz y la salida de audio estaba desactivada: como tal, la entrada sensorial se equiparaba entre el jugador humano y los agentes. El rendimiento humano es la recompensa media obtenida a partir de unos 20 episodios de cada juego con una duración máxima de 5 min cada uno, tras unas 2 h de práctica jugando cada juego.

Algoritmo.Consideraremos tareas en las que un agente interactúa con un entorno, en este caso el emulador de Atari, en una secuencia de acciones, observaciones y recompensas. En cada paso de tiempo el agente selecciona una accióndel conjunto de acciones legales del juego, $A \sim f_1, \dots, f_k$ gramo.La acción se pasa al emulador y modifica su estado interno y la puntuación del juego. En general, el entorno puede ser estocástico. El agente no observa el estado interno del emulador; en cambio, el agente observa una imagen X_t del emulador, que es un vector de valores de píxeles que representan la pantalla actual. Además recibe una recompensaque representa el cambio en la puntuación del juego. Tenga en cuenta que, en general, la puntuación del juego puede depender de toda la secuencia previa de acciones y observaciones; la retroalimentación sobre una acción solo puede recibirse después de que hayan transcurrido muchos miles de pasos de tiempo.

Debido a que el agente solo observa la pantalla actual, la tarea se observa parcialmente³³ y muchos estados del emulador tienen un alias perceptivo (es decir, es imposible comprender completamente la situación actual solo desde la pantalla actual X_t). Por lo tanto, secuencias de acciones y observaciones, $s_t \sim x_1, a_1, x_2, \dots, a_t, x_t$, se ingresan al algoritmo, que luego aprende estrategias de juego dependiendo de estas secuencias. Se supone que todas las secuencias en el emulador terminan en un número finito de pasos de tiempo. Este formalismo da lugar a un proceso de decisión de Markov (MDP) grande pero finito en el que cada secuencia es un estado distinto. Como resultado, podemos aplicar métodos de aprendizaje por refuerzo estándar para MDP, simplemente usando la secuencia completas:

como la representación del estado en ese momento.

El objetivo del agente es interactuar con el emulador seleccionando acciones que maximicen las recompensas futuras. Hacemos la suposición estándar de que las recompensas futuras se descuentan por un factor deCpor paso de tiempo (C se fijó en 0,99 en todo momento), y

X_t^{π} definir el rendimiento descontado futuro en el momento como $R_t = \sum_{t=t}^{\infty} C^{t-t} r_t$, en el cual T es el paso de tiempo en el que termina el juego. Definimos la función acción-valor óptima $q^{\pi}(s, a)$ como el rendimiento máximo esperado que se puede lograr siguiendo cualquier política, después de ver alguna secuencia y luego tomar alguna acción, $q^{\pi}(s, a) = \max_{\pi} \sum_{t=0}^{\infty} C^{t-t} r_t$ — un, π — en el cual π es una política de asignación de secuencias a acciones (o distribuciones sobre acciones).

La función acción-valor óptima obedece a una identidad importante conocida como la ecuación de Bellman. Esto se basa en la siguiente intuición: si el valor óptimo $q^{\pi}(s, a)$ de la secuencia s en el siguiente paso de tiempo fue conocido por todas las acciones posibles a , entonces la estrategia óptima es seleccionar la acción a que maximiza el valor esperado $\sum_a C^{t-t} q^{\pi}(s, a)$:

$$q^{\pi}(s, a) = \max_{a \in A} \sum_a C^{t-t} q^{\pi}(s, a)$$

La idea básica detrás de muchos algoritmos de aprendizaje por refuerzo es estimar la función de valor de acción mediante el uso de la ecuación de Bellman como una actualización iterativa, $q_t^{\pi}(s, a) \leftarrow q_t^{\pi}(s, a) + \alpha [r_t + \gamma \max_{a' \in A} q^{\pi}(s', a') - q_t^{\pi}(s, a)]$. Dichos algoritmos de iteración de valores convergen a la función acción-valor óptima, q^{π} ? q^{π} como q^{π} ? ?. En la práctica, este enfoque básico no es práctico, porque la función acción-valor se estima por separado para cada secuencia, sin ninguna generalización. En cambio, es común usar un aproximador de funciones para estimar la función acción-valor, $Q(s, a) = \sum_a C^{t-t} q^{\pi}(s, a)$. En la comunidad de aprendizaje por refuerzo, este suele ser un aproximador de función lineal, pero

a veces, en su lugar, se utiliza un aproximador de función no lineal, como una red neuronal. Nos referimos a un aproximador de funciones de redes neuronales con pesos como una red Q. La red AQ se puede entrenar ajustando los parámetros en la iteración i para reducir el error cuadrático medio en la ecuación de Bellman, donde el valor óptimo valores objetivos Q^* se sustituyen por valores objetivo aproximados Q tu $\sim rzCmáximo_{a0}q_{s0,a0;h_i}$, usando parámetros Q de alguna iteración anterior.

Esto conduce a una secuencia de funciones de pérdida $L(h)$ que cambia en cada iteración,

$$L(h) \sim \mathbb{E}_{s, a, r \sim p} [y - Q(s, a; h)]^2$$

$$\sim \mathbb{E}_{s, a, r \sim p} [y - Q(s, a; h)]^2 - V_{s0}^2 - \dots$$

Tenga en cuenta que los objetivos dependen de los pesos de la red; esto contrasta con los objetivos utilizados para el aprendizaje supervisado, que se fijan antes de que comience el aprendizaje. En cada etapa de optimización, mantenemos los parámetros de la iteración anterior h_i corregido al optimizar la función de pérdida $L(h)$, resultando en una secuencia de problemas de optimización bien definidos. El término final es la varianza de los objetivos, que no depende de los parámetros. Estamos optimizando actualmente y , por lo tanto, pueden ignorarse. Diferenciando la función de pérdida con respecto a los pesos llegamos al siguiente gradiente:

$$+ h L(h) \sim \mathbb{E}_{s, a, r \sim p} \left[\frac{\partial}{\partial h} Q(s, a; h) \right] - Q(s, a; h) + h Q(s, a; h) : \dots$$

En lugar de calcular las expectativas completas en el gradiente anterior, a menudo es conveniente desde el punto de vista computacional optimizar la función de pérdida mediante el descenso del gradiente estocástico. El conocido algoritmo Q-learning se puede recuperar en este marco actualizando los pesos después de cada paso de tiempo, reemplazando las expectativas usando muestras individuales y estableciendo $i \sim h_i(1)$.

Tenga en cuenta que este algoritmo no tiene modelo: resuelve la tarea de aprendizaje por refuerzo directamente utilizando muestras del emulador, sin estimar explícitamente la recompensa y la dinámica de transición. Puede estar fuera de la política: aprende sobre los codiciosos políticos $\sim argmax_a Q(s, a; h)$, siguiendo una distribución de comportamiento que asegura una adecuada exploración del espacio estatal. En la práctica, la distribución del comportamiento es a menudo seleccionada por una política codiciosa que sigue la política codiciosa con probabilidad 12% y selecciona una acción aleatoria con probabilidad 8%.

Algoritmo de entrenamiento para redes Q profundas. El algoritmo completo para entrenar redes Q profundas se presenta en el Algoritmo 1. El agente selecciona y ejecuta acciones de acuerdo con una política codiciosa basada en q . Debido a que el uso de historias de longitud arbitraria como entradas a una red neuronal puede ser difícil, nuestra función Q funciona en una representación de longitud fija de las historias producidas por la función w descrita arriba. El algoritmo modifica el Q-learning en línea estándar de dos maneras para que sea adecuado para entrenar grandes redes neuronales sin divergir.

Primero, usamos una técnica conocida como repetición de experiencia.²³ En el que almacenamos las experiencias del agente en cada paso de tiempo, $m_i(s_i, a_i, r_i, s_{i+1})$, en un conjunto de datos $D = \{m_1, \dots, m_T\}$, agrupados en muchos episodios (donde el final de un episodio ocurre cuando se alcanza un estado terminal) en una memoria de reproducción. Durante el ciclo interno del algoritmo, aplicamos actualizaciones de Q-learning, o actualizaciones de minibatches, a muestras de experiencia, (s, a, r, s') , extraídas al azar del grupo de muestras almacenadas. Este enfoque tiene varias ventajas sobre el Q-learning en línea estándar. Primero, cada paso de la experiencia se usa potencialmente en muchas actualizaciones de peso, lo que permite una mayor eficiencia de los datos. Segundo, aprender directamente de muestras consecutivas es ineficiente debido a las fuertes correlaciones entre las muestras; la aleatorización de las muestras rompe estas correlaciones y, por lo tanto, reduce la varianza de las actualizaciones. En tercer lugar, cuando se aprende sobre la política, los parámetros actuales determinan la siguiente muestra de datos en la que se entrena los parámetros. Por ejemplo, si la acción de maximización es moverse hacia la izquierda, las muestras de entrenamiento estarán dominadas por muestras del lado izquierdo; si la acción de maximización cambia a la derecha, la distribución del entrenamiento también cambiará.²⁴ Usando experiencia

La distribución del comportamiento se promedia sobre muchos de sus estados anteriores, suavizando el aprendizaje y evitando oscilaciones o divergencias en los parámetros. Tenga en cuenta que cuando se aprende mediante la repetición de la experiencia, es necesario aprender fuera de la política (porque nuestros parámetros actuales son diferentes a los utilizados para generar la muestra), lo que motiva la elección de Q-learning.

En la práctica, nuestro algoritmo solo almacena la última experiencia en la memoria de reproducción y muestra uniformemente al azar de la memoria para realizar actualizaciones. Este enfoque está limitado en algunos aspectos porque el búfer de memoria no diferencia las transiciones importantes y siempre sobrescribe las transiciones recientes debido al tamaño finito de la memoria. De manera similar, el muestreo uniforme otorga la misma importancia a todas las transiciones en la memoria de reproducción. Una estrategia de muestreo más sofisticada podría enfatizar las transiciones de las que podemos aprender más, similar al barrido priorizado.²⁵

La segunda modificación del Q-learning en línea destinada a mejorar aún más la estabilidad de nuestro método con redes neuronales es usar una red separada para generar los objetivos, y , en la actualización de Q-learning. Más precisamente, cada actualización clonamos la red q para obtener una red objetivo Q' que usa Q para generar los objetivos Q-learning y, para el siguiente actualización, Q . Esta modificación hace que el algoritmo sea más estable en comparación con el Q-learning en línea estándar, donde una actualización que aumenta $Q(s, a)$ a menudo también aumenta $Q(s_{t+1}, a)$ para todos y tanto aumenta el objetivo y , lo que posiblemente provoque oscilaciones o divergencias en la política. Generar los objetivos utilizando un conjunto anterior de parámetros agrega un retraso entre el momento en que se actualiza q y el momento en que la actualización afecta a los objetivos y , haciendo que la divergencia o las oscilaciones sean mucho más improbables.

Y-También nos resultó útil recordar el término de error de la actualización $rzCmáximo_{a0}q_{s0,a0;h_i}$ para estar entre 21 y 1. Porque la pérdida de valor absoluto función $|x|$ tiene una derivada de 21 para todos los valores negativos de x y una derivada de 1 para todos los valores positivos de x , recordar el error cuadrático para estar entre 21 y 1 corresponde al uso de una función de pérdida de valor absoluto para errores fuera del (21, 1) intervalo. Esta forma de recorte de errores mejoró aún más la estabilidad del algoritmo. Algoritmo 1: Q-learning profundo con repetición de experiencia. Inicializar la memoria de reproducción a capacidad norte

Inicializar función de valor de acción q con pesos aleatorios h

Inicialice la función de valor de acción objetivo Q con pesos h

Para episodio 51, METRO hacer

Iniciar secuencias $\sim f(x)$ gramoy secuencia preprocesada $w \sim w - s - b$ Para

t51, Thacher

con probabilidad de elegir una acción aleatoria a :

de lo contrario seleccionar $\sim argmax_a Q(w - s - b, a; h)$

Ejecutar acción a en el emulador y observar la recompensa r e imagen x_{t+1}

Establecer $s_{t+1} \sim s_t, a_t, x_{t+1}$ y preprocesar $w_{t+1} \sim w - s_{t+1} - b$

Transición de la tienda w, a_t, r, t_{t+1} y w_{t+1} $\sim \sim \sim$

Ejemplo de minibatch aleatorio de transiciones w_j, a_j, r_j, t_{j+1} de D

($j \sim \sim \sim$) si el episodio termina en el paso j

Establecer $y_j = Q(w_j, a_j; h)$ de lo contrario

$r_j - rzCmáximo_{a0}Q(w_j, a_j; h)$ $\sim \sim \sim$

Realice un paso de descenso de gradiente en y_j ($w_j, a_j; h$) $\sim \sim \sim$ Con respecto a

parámetros de red h

Cada $Capos$ restablecer $Q = q$

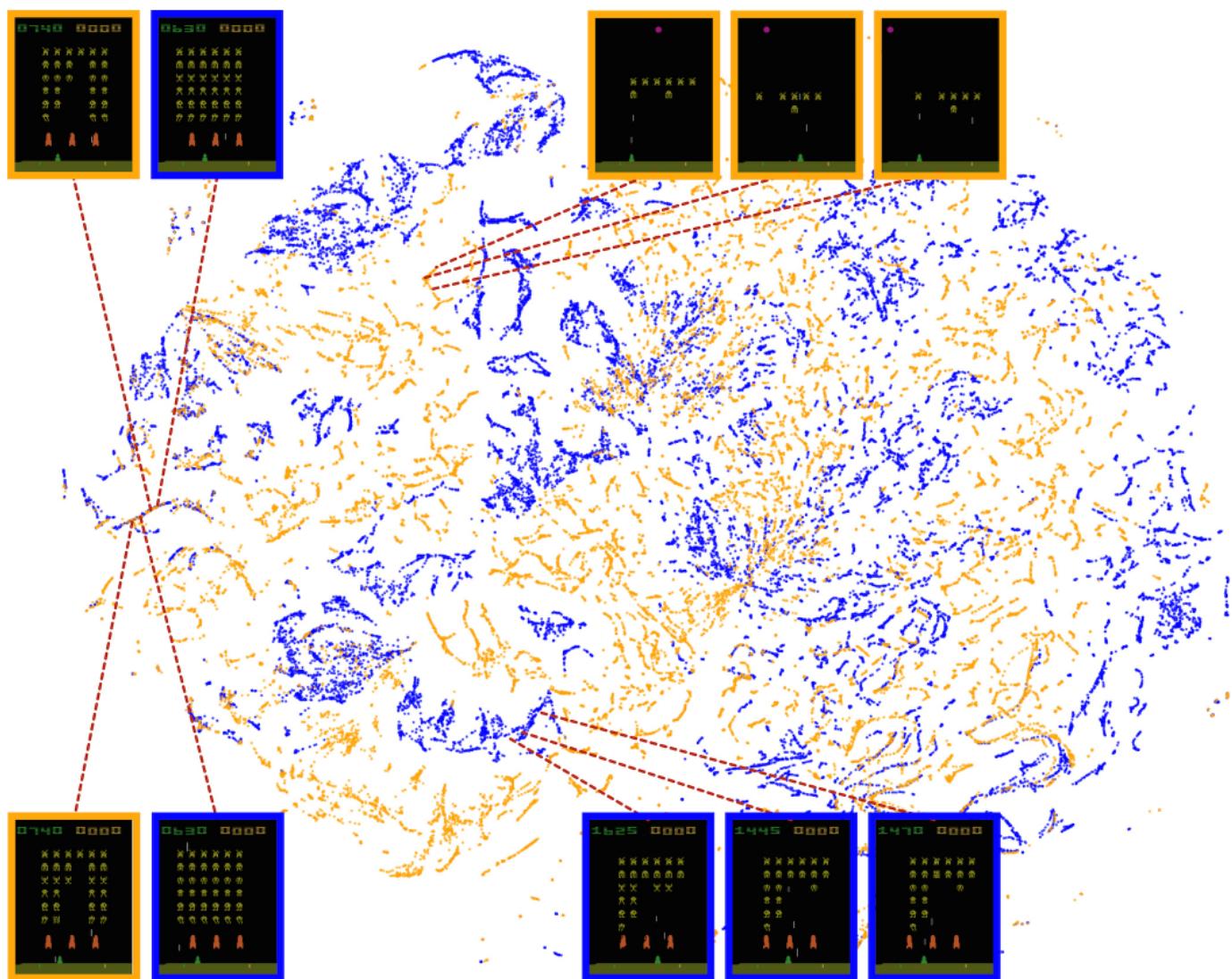
Terminar por

Terminar por

31. Jarrett, K., Kavukcuoglu, K., Ranzato, M.A. y LeCun, Y. ¿Cuál es la mejor arquitectura de varias etapas para el reconocimiento de objetos? Proc. IEEE. Int. Conf. Comput. Vis. 2146-2153 (2009).

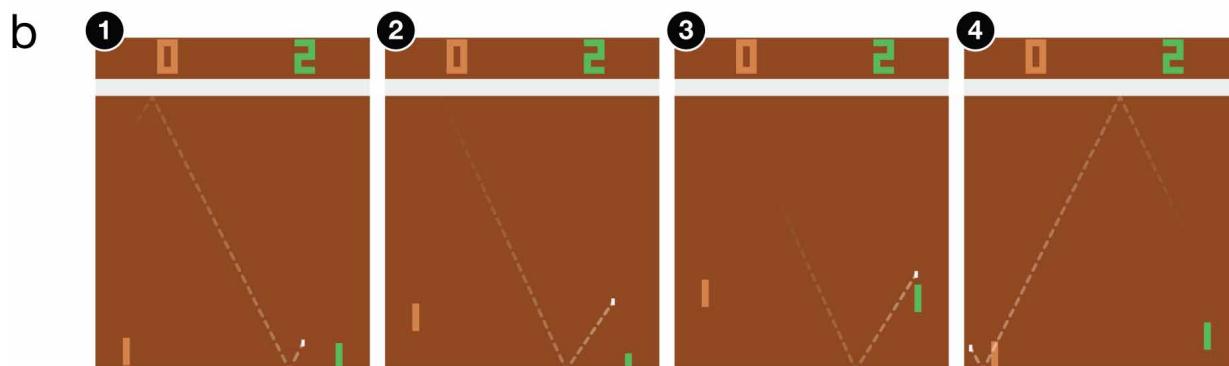
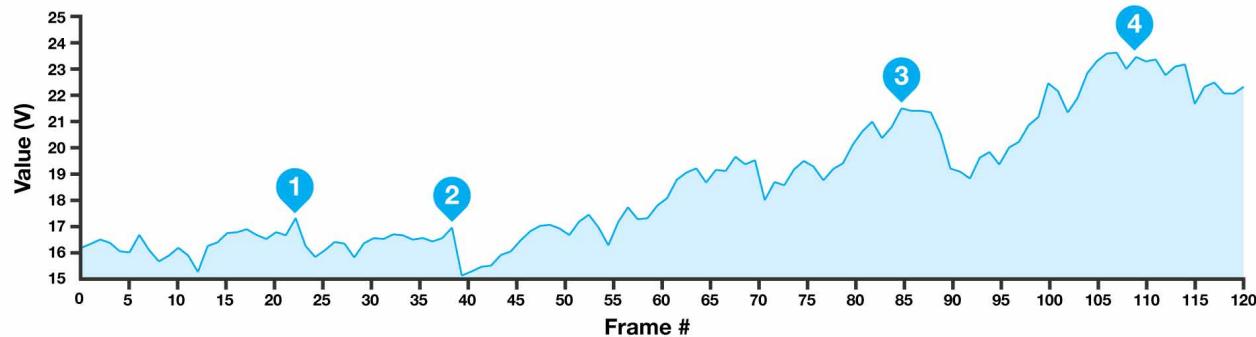
32. Nair, V. & Hinton, G.E. Las unidades lineales rectificadas mejoran las máquinas Boltzmann restringidas. Proc. En t. Conf. Mach. Aprender. 807-814 (2010).

33. Kaelbling, L.P., Littman, M.L. y Cassandra, A.R. Planificación y actuación en dominios estocásticos parcialmente observables. Inteligencia artificial 101, 99-134 (1994).



Datos extendidos Figura 1 | Incorporación bidimensional de t-SNE de las representaciones en la última capa oculta asignada por DQN a los estados del juego experimentados durante una combinación de juego humano y de agente en Space Invaders. El gráfico se generó ejecutando el algoritmo t-SNE₂₅ en la representación de la última capa oculta asignada por DQN a los estados del juego experimentados durante una combinación de juego humano (30 min) y agente (2 h). El hecho de que haya una estructura similar en las incrustaciones bidimensionales correspondientes a la representación DQN de estados experimentados durante el juego humano (naranja

puntos) y el juego DQN (puntos azules) sugiere que las representaciones aprendidas por DQN sí se generalizan a los datos generados a partir de políticas distintas a las suyas. La presencia en la incrustación de t-SNE de grupos superpuestos de puntos correspondientes a la representación de la red de estados experimentados durante el juego humano y del agente muestra que el agente DQN también sigue secuencias de estados similares a las que se encuentran en el juego humano. Se muestran capturas de pantalla correspondientes a los estados seleccionados (humano: borde naranja; DQN: borde azul).



Datos extendidos Figura 2 | Visualización de funciones de valor aprendido en dos juegos, Breakout y Pong. a, Una visualización de la función de valor aprendido en el juego Breakout. En los puntos de tiempo 1 y 2, se prevé que el valor de estado sea 17 y el agente está limpiando los ladrillos en el nivel más bajo. Cada uno de los picos en la curva de función de valor corresponde a una recompensa obtenida al despejar un ladrillo. En el punto de tiempo 3, el agente está a punto de atravesar el nivel superior de ladrillos y el valor aumenta a .21 en previsión de romper y limpiar un gran conjunto de ladrillos. En el punto 4, el valor está por encima de 23 y el agente se ha abierto paso. Pasado este punto, la pelota rebota en la parte superior de los ladrillos, despejando muchos de ellos por sí sola. b, Una visualización de la función acción-valor aprendida en el juego Pong. En el punto de tiempo 1, la pelota se mueve hacia la paleta controlada por el agente en el lado derecho de la pantalla y los valores de

todas las acciones están alrededor de 0.7, lo que refleja el valor esperado de este estado basado en la experiencia previa. En el punto de tiempo 2, el agente comienza a mover la paleta hacia la pelota y el valor de la acción "arriba" se mantiene alto mientras que el valor de la acción "abajo" cae a 0.9. Esto refleja el hecho de que presionar 'abajo' haría que el agente perdiera el balón e incurriera en una recompensa de -21. En el punto de tiempo 3, el agente golpea la pelota presionando 'arriba' y la recompensa esperada sigue aumentando hasta el punto de tiempo 4, cuando la pelota llega al borde izquierdo de la pantalla y el valor de todas las acciones refleja que el agente está a punto para recibir una recompensa de 1. Tenga en cuenta que la línea discontinua muestra la trayectoria pasada de la pelota únicamente con fines ilustrativos (es decir, no se muestra durante el juego). Con permiso de Atari Interactive, Inc.

Tabla de datos extendida 1 | Lista de hiperparámetros y sus valores

Hyperparameter	Value	Description
minibatch size	32	Number of training cases over which each stochastic gradient descent (SGD) update is computed.
replay memory size	1000000	SGD updates are sampled from this number of most recent frames.
agent history length	4	The number of most recent frames experienced by the agent that are given as input to the Q network.
target network update frequency	10000	The frequency (measured in the number of parameter updates) with which the target network is updated (this corresponds to the parameter C from Algorithm 1).
discount factor	0.99	Discount factor gamma used in the Q-learning update.
action repeat	4	Repeat each action selected by the agent this many times. Using a value of 4 results in the agent seeing only every 4th input frame.
update frequency	4	The number of actions selected by the agent between successive SGD updates. Using a value of 4 results in the agent selecting 4 actions between each pair of successive updates.
learning rate	0.00025	The learning rate used by RMSProp.
gradient momentum	0.95	Gradient momentum used by RMSProp.
squared gradient momentum	0.95	Squared gradient (denominator) momentum used by RMSProp.
min squared gradient	0.01	Constant added to the squared gradient in the denominator of the RMSProp update.
initial exploration	1	Initial value of ϵ in ϵ -greedy exploration.
final exploration	0.1	Final value of ϵ in ϵ -greedy exploration.
final exploration frame	1000000	The number of frames over which the initial value of ϵ is linearly annealed to its final value.
replay start size	50000	A uniform random policy is run for this number of frames before learning starts and the resulting experience is used to populate the replay memory.
no-op max	30	Maximum number of "do nothing" actions to be performed by the agent at the start of an episode.

Los valores de todos los hiperparámetros se seleccionaron realizando una búsqueda informal en los juegos Pong, Breakout, Seaquest, Space Invaders y BeamRider. No realizamos una búsqueda de cuadrícula sistemática debido al alto costo computacional, aunque es concebible que se puedan obtener resultados aún mejores ajustando sistemáticamente los valores de los hiperparámetros.

Tabla de datos extendida 2 | Comparación de puntuajes de juegos obtenidos por agentes DQN con métodos de la literatura^{12,15} y un probador profesional de juegos humanos

Game	Random Play	Best Linear Learner	Contingency (SARSA)	Human	DQN (\pm std)	Normalized DQN (% Human)
Alien	227.8	939.2	103.2	6875	3069 (\pm 1093)	42.7%
Amidar	5.8	103.4	183.6	1676	739.5 (\pm 3024)	43.9%
Assault	222.4	628	537	1496	3359 (\pm 775)	246.2%
Asterix	210	987.3	1332	8503	6012 (\pm 1744)	70.0%
Asteroids	719.1	907.3	89	13157	1629 (\pm 542)	7.3%
Atlantis	12850	62687	852.9	29028	85641 (\pm 17600)	449.9%
Bank Heist	14.2	190.8	67.4	734.4	429.7 (\pm 650)	57.7%
Battle Zone	2360	15820	16.2	37800	26300 (\pm 7725)	67.6%
Beam Rider	363.9	929.4	1743	5775	6846 (\pm 1619)	119.8%
Bowling	23.1	43.9	36.4	154.8	42.4 (\pm 88)	14.7%
Boxing	0.1	44	9.8	4.3	71.8 (\pm 8.4)	1707.9%
Breakout	1.7	5.2	6.1	31.8	401.2 (\pm 26.9)	1327.2%
Centipede	2091	8803	4647	11963	8309 (\pm 5237)	63.0%
Chopper Command	811	1582	16.9	9882	6687 (\pm 2916)	64.8%
Crazy Climber	10781	23411	149.8	35411	114103 (\pm 22797)	419.5%
Demon Attack	152.1	520.5	0	3401	9711 (\pm 2406)	294.2%
Double Dunk	-18.6	-13.1	-16	-15.5	-18.1 (\pm 2.6)	17.1%
Enduro	0	129.1	159.4	309.6	301.8 (\pm 24.6)	97.5%
Fishing Derby	-91.7	-89.5	-85.1	5.5	-0.8 (\pm 19.0)	93.5%
Freeway	0	19.1	19.7	29.6	30.3 (\pm 0.7)	102.4%
Frostbite	65.2	216.9	180.9	4335	328.3 (\pm 250.5)	6.2%
Gopher	257.6	1288	2368	2321	8520 (\pm 3279)	400.4%
Gravitar	173	387.7	429	2672	306.7 (\pm 223.9)	5.3%
H.E.R.O.	1027	6459	7295	25763	19950 (\pm 158)	76.5%
Ice Hockey	-11.2	-9.5	-3.2	0.9	-1.6 (\pm 2.5)	79.3%
James Bond	29	202.8	354.1	406.7	576.7 (\pm 175.5)	145.0%
Kangaroo	52	1622	8.8	3035	6740 (\pm 2959)	224.2%
Krull	1598	3372	3341	2395	3805 (\pm 1033)	277.0%
Kung-Fu Master	258.5	19544	29151	22736	23270 (\pm 5955)	102.4%
Montezuma's Revenge	0	10.7	259	4367	0 (\pm 0)	0.0%
Ms. Pacman	307.3	1692	1227	15693	2311 (\pm 525)	13.0%
Name This Game	2292	2500	2247	4076	7257 (\pm 547)	278.3%
Pong	-20.7	-19	-17.4	9.3	18.9 (\pm 1.3)	132.0%
Private Eye	24.9	684.3	86	69571	1788 (\pm 5473)	2.5%
Q*Bert	163.9	613.5	960.3	13455	10596 (\pm 3294)	78.5%
River Raid	1339	1904	2650	13513	8316 (\pm 1049)	57.3%
Road Runner	11.5	67.7	89.1	7845	18257 (\pm 4268)	232.9%
Robotank	2.2	28.7	12.4	11.9	51.6 (\pm 4.7)	509.0%
Seaquest	68.4	664.8	675.5	20182	5286 (\pm 1310)	25.9%
Space Invaders	148	250.1	267.9	1652	1976 (\pm 893)	121.5%
Star Gunner	664	1070	9.4	10250	57997 (\pm 3152)	598.1%
Tennis	-23.8	-0.1	0	-8.9	-2.5 (\pm 1.9)	143.2%
Time Pilot	3568	3741	24.9	5925	5947 (\pm 1600)	100.9%
Tutankham	11.4	114.3	98.2	167.6	186.7 (\pm 41.9)	112.2%
Up and Down	533.4	3533	2449	9082	8456 (\pm 3162)	92.7%
Venture	0	66	0.6	1188	380.0 (\pm 238.6)	32.0%
Video Pinball	16257	16871	19761	17298	42684 (\pm 16287)	2539.4%
Wizard of Wor	563.5	1981	36.9	4757	3393 (\pm 2019)	67.5%
Zaxxon	32.5	3365	21.4	9173	4977 (\pm 1235)	54.1%

Best Linear Learner es el mejor resultado obtenido por un aproximador de funciones lineales en diferentes tipos de características diseñadas a mano. Las cifras de agentes de contingencia (SARSA) son los resultados obtenidos en la ref. 15. Tenga en cuenta que las cifras de la última columna indican el rendimiento de DQN en relación con el probador de juegos humanos, expresado como porcentaje, es decir, $1003 / (\text{puntuación DQN} + \text{puntuación de juego aleatorio}) / (\text{puntuación humana} + \text{puntuación de juego aleatorio})$.

Tabla de datos extendida 3 | Los efectos de la reproducción y la separación de la red Q de destino

Game	With replay, with target Q	With replay, without target Q	Without replay, with target Q	Without replay, without target Q
Breakout	316.8	240.7	10.2	3.2
Enduro	1006.3	831.4	141.9	29.1
River Raid	7446.6	4102.8	2867.7	1453.0
Seaquest	2894.4	822.6	1003.0	275.8
Space Invaders	1088.9	826.3	373.2	302.0

Se capacitó a los agentes de DQN para 10 millones de fotogramas utilizando hiperparámetros estándar para todas las combinaciones posibles de activación o desactivación de la reproducción, uso o no uso de una red Q de destino separada y tres tasas de aprendizaje diferentes. Cada agente se evaluó cada 250 000 cuadros de entrenamiento para 135 000 cuadros de validación y se informa el puntaje promedio más alto del episodio. Tenga en cuenta que estos episodios de evaluación no se truncaron a los 5 minutos, lo que llevó a puntuaciones más altas en Enduro que las que se informan en la Tabla 2 de datos ampliados. Tenga en cuenta también que el número de fotogramas de entrenamiento fue más corto (10 millones de fotogramas) en comparación con los resultados principales presentados en los datos ampliados. Tabla 2 (50 millones de fotogramas).

Tabla de datos extendida 4 | Comparación del rendimiento de DQN con el aproximador de función lineal

Game	DQN	Linear
Breakout	316.8	3.00
Enduro	1006.3	62.0
River Raid	7446.6	2346.9
Seaquest	2894.4	656.9
Space Invaders	1088.9	301.3

El rendimiento de DQNagent se compara con el rendimiento de un aproximador de función lineal en los 5 juegos de validación (es decir, donde se usó una sola capa lineal en lugar de la red convolucional, en combinación con reproducción y red objetivo separada). Se capacitó a los agentes para 10 millones de fotogramas utilizando hiperparámetros estándar y tres tasas de aprendizaje diferentes. Cada agente se evaluó cada 250 000 cuadros de entrenamiento para 135 000 cuadros de validación y se informa el puntaje promedio más alto del episodio. Tenga en cuenta que estos episodios de evaluación no se truncaron a los 5 minutos, lo que llevó a puntuaciones más altas en Enduro que las que se informan en la Tabla 2 de datos ampliados. Tenga en cuenta también que el número de fotogramas de entrenamiento fue más corto (10 millones de fotogramas) en comparación con los resultados principales presentados en los datos ampliados. Tabla 2 (50 millones de fotogramas).