

E

Revista Brasileira de

Engenharias

ISSN 3085-8089

vol. 1, n. 2, 2025

... ARTIGO 5

Data de Aceite: 10/12/2025

APLICAÇÃO DE APRENDIZAGEM POR REFORÇO ADAPTATIVO PARA O DESVIO DE OBSTÁCULOS EM AGENTES VIRTUAIS

Márcio Mendonça

Universidade Tecnológica Federal do Paraná
PPGEM-CP - Programa de Pós-Graduação em Engenharia Mecânica CP/PG
Cornélio Procópio - PR
<http://lattes.cnpq.br/5415046018018708>

Vitor Blanc Milani

Universidade Tecnológica Federal do Paraná
Mestrando - PPGEM-CP - Programa de Pós-Graduação em Engenharia Mecânica CP/PG
Cornélio Procópio - PR
<http://lattes.cnpq.br/4504374098250296>

Juliana Maria de Jesus Ribeiro

Universidade Tecnológica Federal do Paraná (UTFPR) – Campus Cornélio Procópio/Londrina, Paraná
– Brasil
Mestranda no Programa de Pós-Graduação em Ensino de Ciências Humanas, Sociais e da Natureza
– PPGEN
Londrina - PR
<http://lattes.cnpq.br/6279504657014354>

Fabio Rodrigo Milanez

UnISENAIPR-Campus Londrina
Londrina-PR
<http://lattes.cnpq.br/3808981195212391>

Marcos Dantas de Oliveira

CEAL- Clube de Engenharia de Londrina
Londrina-PR
<https://lattes.cnpq.br/5329306535174160>



Todo o conteúdo desta revista está licenciado sob a Licença Creative Commons Atribuição 4.0 Internacional (CC BY 4.0).

Iago Maran Machado

Universidade Tecnológica Federal do Paraná
Mestrando (aluno externo) - PPGEM-CP -
Programa de Pós-Graduação em Engenharia
Mecânica CP/PG
Cornélio Procópio - PR
<http://lattes.cnpq.br/3808981195212391>

Emerson Ravazzi Pires da Silva

Universidade Tecnológica Federal do Paraná
Departamento Acadêmico de Engenharia
Elétrica (DAELE)
Cornélio Procópio - PR
<http://lattes.cnpq.br/3845751794448092>

Eduardo Pegoraro Heinemann

Universidade Tecnológica Federal do Paraná
Departamento Acadêmico de Engenharia
Elétrica (DAELE)
Cornélio Procópio - PR
<http://lattes.cnpq.br/0964474292409084>

Marco Antônio Ferreira Finocchio

Universidade Tecnológica Federal do Paraná
Departamento Acadêmico de Engenharia
Elétrica (DAELE)
Cornélio Procópio - PR
<http://lattes.cnpq.br/8619727190271505>

Andressa Haiduk

Dimension Engenharia
Rio Negro - PR
<http://lattes.cnpq.br/2786786167224165>

Junior Candido Mendonça

Discente-Universidade Tecnológica Federal
do Paraná Departamento
Acadêmico de Engenharia Elétrica (DAELE)
Cornélio Procópio - PR
<http://lattes.cnpq.br/9637563407033947>

Francisco de Assis Scannavino Junior

Universidade Tecnológica Federal do Paraná
Departamento Acadêmico de Engenharia
Elétrica (DAELE)
Cornélio Procópio - PR
<http://lattes.cnpq.br/4513330681918118>

Tatiane Monteiro Pereira

Mestranda - Programa de Pós-Graduação em
Ensino de Ciências Humanas, Sociais e da
Natureza (PPGEN-CP/LD)
Cornélio Procópio - PR
<http://lattes.cnpq.br/9520601026438758>

Vera Adriana Huang Azevedo

Hypólito

Centro Estadual de Educação Tecnológica
Paula Souza
Etec Jacinto Ferreira de Sá
Ourinhos - SP
<http://lattes.cnpq.br/6169590836932698>

Angelo Feracin Neto

Universidade Tecnológica Federal do Paraná
Departamento Acadêmico de Engenharia
Elétrica (DAELE)
Cornélio Procópio - PR
<http://lattes.cnpq.br/0580089660443472>

Roberto Bondarik

Programa de Pós-Graduação em Ensino de
Ciências Humanas, Sociais e da Natureza
(PPGEN-CP/LD)
Cornélio Procópio - PR
<http://lattes.cnpq.br/6263028023417758>

André Luiz Salvat Moscato

Instituto Federal do Paraná, Campus Jacarez-
inho
Jacarezinho - PR
<http://lattes.cnpq.br/1744149363927228>

Ricardo Breganon

Instituto Federal do Paraná, Campus Jacarez-
inho
Jacarezinho - PR
<http://lattes.cnpq.br/2441043775335349>

Fabio Nogueira de Queiroz

Centro Paula Souza
Departamento Computação-FATEC
<http://lattes.cnpq.br/4466493001956276>

Armando Paulo da Silva

Programa de Pós-Graduação em Ensino de
Ciências Humanas, Sociais e da Natureza
(PPGEN-CP/LD)
Cornélio Procópio - PR
<http://lattes.cnpq.br/6724994186659242>

Mário Sérgio Martinelli Medina

Técnico de Laboratório de Informática
<http://lattes.cnpq.br/6624259960900455>

Resumo: O aprendizado de máquina constitui um dos principais eixos da Inteligência Artificial contemporânea, ao permitir que sistemas computacionais realizem inferências e tomem decisões com base em dados. Entre suas diferentes abordagens — aprendizado supervisionado, não supervisionado e por reforço — destaca-se a aprendizagem por reforço, que utiliza o paradigma agente-ambiente para otimizar políticas de ação por meio de recompensas acumuladas ao longo do tempo. Nesse contexto, a aprendizagem por reforço adaptativa exerce um papel essencial, pois capacita o agente a ajustar seu comportamento em ambientes dinâmicos, não estacionários ou sujeitos a incertezas. Ao incorporar mecanismos adaptativos, o agente deixa de operar com uma política estática e passa a modificar suas estratégias de forma contínua, ampliando sua capacidade de resposta diante de mudanças estruturais e perturbações do ambiente. Essa característica torna o método especialmente adequado para aplicações em robótica móvel, controle inteligente, veículos autônomos, manufatura avançada, jogos e processos industriais complexos. Além disso, técnicas adaptativas em reforço incrementam a robustez, a capacidade de generalização e a resiliência dos modelos. Estratégias como ajuste dinâmico entre exploração e exploração, meta-aprendizagem, otimização em tempo real e integração com aprendizado profundo elevam o desempenho e a autonomia dos agentes. Conclui-se que a aprendizagem por reforço adaptativa representa um componente crucial para o desenvolvimento de sistemas inteligentes capazes de operar com eficiência, segurança e estabilidade em ambientes complexos e variáveis.

Palavras-Chave: Aprendizado de máquina. Aprendizagem por reforço. Sistemas adap-

tativos. Inteligência Artificial. Autonomia computacional.

INTRODUÇÃO

A educação do século XXI demanda abordagens inovadoras que superem o ensino fragmentado em disciplinas estáticas. O modelo tradicional, ainda presente em muitas instituições, mostra-se insuficiente para atender às exigências de uma sociedade marcada pela complexidade, pela circulação em rede de informações e pela velocidade das transformações tecnológicas (MORAN, 2018).

Nesse contexto, a aprendizagem por reforço (*Reinforcement Learning – RL*) consolidou-se, nas últimas décadas, como uma das metodologias mais relevantes dentro do campo da Inteligência Artificial (IA), como por exemplo o *Deepseek*, I.A. Generativa (DEEPSEEK, 2025) para o desenvolvimento de agentes autônomos capazes de aprender por interação direta com o ambiente. Diferentemente dos paradigmas supervisionado e não supervisionado, o RL opera por meio da maximização de recompensas acumuladas ao longo do tempo, estruturando o processo de aprendizado a partir da relação dinâmica agente-ambiente (SUTTON; BARTO, 2018). Essa característica concede à abordagem elevada eficiência em tarefas que envolvem tomada de decisão sequencial sob incerteza, como navegação robótica, controle inteligente, jogos, sistemas ciber físicos e veículos autônomos.

Apesar de seu potencial, a aplicação de RL a cenários reais evidencia desafios importantes relacionados à não-estacionalidade do ambiente, à alta dimensionalidade dos estados e à necessidade de adaptação contínua das políticas de controle. Ambientes

dinâmicos, a exemplo de sistemas robóticos em operação, apresentam variações estruturais que exigem que o agente seja capaz de ajustar sua política em tempo real, sob pena de perda de desempenho ou até falha operacional. A literatura demonstra que políticas estáticas, treinadas *off-line*, podem apresentar degradação significativa de desempenho quando expostas a perturbações previamente não observadas (HASSANZADEH; ZHANG; HA, 2022; ZHOU et al., 2020). Diante desses limites, emerge a aprendizagem por reforço adaptativa, que amplia o paradigma clássico ao permitir que o agente atualize parâmetros, redefine estratégias de exploração e revise modelos internos à medida que o ambiente se transforma (NAGABANDI et al., 2018).

A aprendizagem por reforço adaptativa incorpora princípios do RL tradicional, combinando-os com técnicas de meta-aprendizagem, ajustes dinâmicos de políticas e mecanismos de atualização contínua, o que resulta em maior robustez diante de incertezas e variabilidades ambientais. Estratégias adaptativas como o ajuste automático da taxa exploração/exploração (ϵ -greedy dinâmico ou UCB adaptativo), políticas parametrizadas que se reconfiguram diante de novos estados, ou ainda mecanismos de aprendizagem hierárquica permitem ao agente refinar seu comportamento mesmo após o treinamento inicial (KIRK et al., 2023). Esses avanços tornam o método adequado para aplicações em que as condições ambientais mudam rapidamente, como navegação de robôs móveis, controle de drones, sistemas de manufatura flexível e robôs autônomos em ambientes não estruturados.

No domínio da robótica e dos agentes virtuais, o desvio de obstáculos constitui uma das tarefas mais investigadas e essen-

ciais para garantir segurança e autonomia. Métodos clássicos, como campos potenciais ou planejamento geométrico, embora eficientes em cenários simples, apresentam limitações quando expostos a obstáculos dinâmicos, comportamentos imprevisíveis ou ambientes parcialmente observáveis (SILVA; COSTA; ROSSI, 2021). Nesse panorama, abordagens baseadas em RL adaptativo oferecem vantagens expressivas, ao permitir que o agente construa, refine e ajuste sua política de navegação a partir de experiências sucessivas, desenvolvendo não apenas a capacidade de evitar colisões, mas também de antecipar situações de risco, otimizar trajetórias e adaptar-se a perturbações externas.

O texto-base desta pesquisa (Aplicação de Aprendizagem por Reforço) destaca a importância da adaptação contínua para melhorar o desempenho de agentes virtuais em tarefas de desvio de obstáculos, ressaltando que políticas estáticas frequentemente se mostram insuficientes quando expostas a mudanças estruturais no ambiente. A literatura converge nessa perspectiva ao demonstrar que mecanismos adaptativos elevam significativamente a capacidade de generalização dos agentes, reduzindo o risco de *overfitting* às condições de treinamento e ampliando a resiliência do sistema (MOHAMMED; BEER; WANG, 2022). Além disso, a integração de RL adaptativo com redes neurais profundas (*Deep RL*) tem impulsionado resultados de alto desempenho em domínios complexos, como navegação tridimensional, ambientes com obstáculos móveis e simulações realistas baseadas em física (MNIH et al., 2015; KALASHNIKOV et al., 2018).

Assim, ao investigar a aplicação de aprendizagem por reforço adaptativa para o desvio de obstáculos em agentes virtuais,

esta pesquisa contribui para o desenvolvimento de soluções inteligentes mais robustas, eficientes e escaláveis, especialmente em cenários marcados por alta imprevisibilidade e complexidade ambiental, nos quais agentes autônomos necessitam de comportamento dinâmico e responsivo. A evolução dessa abordagem pode impactar de forma direta áreas como robótica educacional, robótica industrial, veículos autônomos, simulações científicas e sistemas ciberfísicos, ampliando seu potencial de aplicação práticas e aprimorando o estado da arte em agentes inteligentes.

Este artigo organiza-se da seguinte forma: a seção 2 apresenta a revisão bibliográfica do tema; a seção 3 expõe a fundamentação teórica; a seção 4 apresenta e analisa os resultados; e a seção 5 reúne as conclusões e indica possibilidades para trabalhos futuros.

REVISÃO BIBLIOGRÁFICA

A aprendizagem por reforço (*Reinforcement Learning* – RL) consolidou-se como um dos pilares da Inteligência Artificial (IA) moderna, sendo aplicada em domínios como robótica móvel, jogos, sistemas autônomos e otimização industrial. Fundamentada no paradigma agente-ambiente, a RL busca determinar políticas capazes de maximizar recompensas esperadas ao longo do tempo, configurando-se como um processo iterativo baseado em tentativa e erro. Nas últimas décadas, os modelos evoluíram de algoritmos tabulares simples para arquiteturas profundas de grande escala, sobretudo após o surgimento da aprendizagem por reforço profundo (*Deep Reinforcement Learning* – DRL), que combina RL com redes neurais artificiais para permitir generalização em es-

paços de estados contínuos e de alta dimensionalidade (MNIH et al., 2015; LI, 2017).

1. Fundamentos da Aprendizagem por Reforço

A RL é tradicionalmente estruturada por quatro componentes fundamentais: o agente, o ambiente, a política de decisão e a função de recompensa. A interação entre esses elementos é modelada, em geral, como um Processo de Decisão de Markov (MDP). Entre os algoritmos clássicos, destacam-se *Q-Learning* (WATKINS; DAYAN, 1992), SARSA (RUMMERY; NIRANJAN, 1994) e *Actor-Critic* (KONIDARIS; BARTO, 2007), amplamente utilizados como base para aplicações complexas.

Contudo, em ambientes dinâmicos ou não estacionários — como missões de robôs móveis, navegação autônoma ou desvio de obstáculos — tais algoritmos tradicionais são insuficientes, dado que dependem de hipóteses de estacionariedade e estabilidade do ambiente. Isso motivou o desenvolvimento de abordagens adaptativas capazes de ajustar parâmetros de política, taxas de aprendizado e estratégias de exploração ao longo do processo (HAN; WANG, 2020).

2. Aprendizagem por Reforço Adaptativa

A aprendizagem por reforço adaptativa (*Adaptive Reinforcement Learning* – ARL) surge como uma alternativa ao permitir que o agente modifique sua política em tempo real, com base nas mudanças ambientais, perturbações estruturais ou novas condições não previamente observadas. Tal capacidade é essencial para cenários reais, como robótica autônoma, veículos autoguiados e sistemas inteligentes de manufatura (KIUMARSI et al., 2020).

Entre as técnicas adaptativas mais utilizadas, destacam-se:

Ajuste adaptativo entre exploração

Métodos que equilibram automaticamente exploração (buscar novos comportamentos) e exploração (usar comportamentos aprendidos) são particularmente importantes em ambientes com mudanças rápidas. Estratégias como *ϵ -greedy decay*, *Upper Confidence Bounds* (UCB) e *Softmax adaptativo* são utilizadas para regular esse equilíbrio dinamicamente (AUER, 2002; TOKIC, 2010).

Meta-aprendizagem

Modelos de meta-aprendizagem permitem que o agente aprenda a aprender (seja condicionado a aprender), atualizando rapidamente seus parâmetros internos com poucas experiências. Finn, Abbeel e Levine (2017) demonstram essa capacidade por meio do algoritmo *Model-Agnostic Meta-Learning* (MAML), que se tornou referência na área.

Adaptabilidade baseada em gradientes

A incorporação de redes profundas permite que agentes ajustem não apenas políticas, mas representações internas, favorecendo a adaptação a novos domínios (ZHANG et al., 2021). Essa abordagem tem sido aplicada com sucesso em políticas robustas para locomoção e navegação.

Controle adaptativo associado à RL

Em robótica, métodos híbridos que combinam RL com controle adaptativo tradicional (como controle robusto ou PID

adaptativo) têm apresentado resultados promissores, principalmente para garantir estabilidade e segurança (KUPCSIK et al., 2013; MODARES; LEWIS, 2014).

Aprendizagem por Reforço em Navegação Autônoma e Desvio de Obstáculos

No domínio da navegação autônoma e no desvio de obstáculos a RL consolidou-se como uma das técnicas mais promissora, especialmente em robôs móveis, drones e veículos autônomos. Em ambientes com obstáculos, a adaptação contínua é crucial, pois a distribuição dos estados muda de forma não determinística. Works como Zhi et al. (2019) e Tai et al. (2017) demonstram que agentes treinados com DRL podem aprender comportamentos complexos, como desvio de obstáculos, mesmo em cenários parcialmente observáveis.

Outra linha relevante envolve as abordagens simuladas-para-real (*sim-to-real*) têm sido amplamente estudadas (JAMES; DAVIDSON; JOHNSON, 2019), permitindo que agentes adaptativos transferem conhecimento de ambientes simulados para cenários reais, ajustando políticas conforme diferenças estruturais emergem.

Pesquisas recentes indicam que agentes com mecanismos adaptativos possuem maior resiliência, estabilidade e robustez frente a mudanças abruptas (PATTANAIK et al., 2018). Em contrapartida, tais sistemas apresentam desafios, incluindo: maior custo computacional, dificuldade em garantir estabilidade global; maior variância durante a fase de exploração adaptativa e riscos de sobre ajuste em ambientes específicos.

Ainda assim, evidências experimentais demonstram que a RL adaptativa supera significativamente métodos estáticos quando aplicada em ambientes complexos e dinâmicos, como tráfego urbano, controle de manipuladores robóticos e navegação em florestas de obstáculos (KIM; KIM; OH, 2022).

O texto fornecido no arquivo destaca a relevância da aprendizagem por reforço adaptativa para agentes em ambientes dinâmicos e incertos, especialmente em tarefas de desvio de obstáculos. Essa perspectiva alinha-se ao consenso atual da área, que considera a adaptação contínua um elemento central para garantir autonomia real (*real-world autonomy*), segurança e capacidade de generalização de sistemas robóticos e inteligentes.

Também se observa a convergência com pesquisas que defendem a necessidade de adaptação contínua para garantir segurança, eficiência e capacidade de generalização aspectos cada vez mais centrais no desenvolvimento de IA moderna.

Aprendizagem por Reforço Adaptativa

A Aprendizagem por Reforço (*Reinforcement Learning – RL*) constitui um paradigma de aprendizado em que um agente aprende a tomar decisões sequenciais a partir da interação com um ambiente incerto, recebendo recompensas e ajustando sua política de ação para maximizar o retorno acumulado ao longo do tempo. Diferentemente do aprendizado supervisionado, a RL não há pares entrada-saída rotulados; o conhecimento é adquirido por tentativa e erro, por meio do feedback escalar de re-

compensa para aprimorar progressivamente seu comportamento.

Formalmente, o problema é modelado como um **Processo de Decisão de Markov (MDP)**, definido pelo quintuplo, em que é o conjunto de estados, o conjunto de ações, a dinâmica estocástica do ambiente, a recompensa imediata e o fator de desconto aplicado às recompensas futuras. O objetivo é encontrar uma política ótima que maximize o **retorno esperado**:

Fundamentação Completa da Aprendizagem por Reforço com Equações Explicadas

Retorno Esperado

O retorno esperado representa a soma descontada de todas as recompensas futuras.

Equação:

$$G_t = \sum \gamma^k R_{t+k+1}$$

Essa soma ponderada define o quanto uma trajetória futura vale para o agente, considerando que recompensas distantes devem valer menos que recompensas imediatas.

Função de Valor de Estado

A função de valor indica a qualidade de estar em um estado s .

Equação:

$$v\pi(s) = E\pi[G_t | S_t = s]$$

Ela mede o valor médio esperado futuro ao seguir a política π a partir de s .

Equação de Bellman

A equação de Bellman define a relação recursiva entre valores:

$$v\pi(s) = \sum \pi(a|s) \sum p(s',r|s,a) [r + \gamma v\pi(s')]$$

Isso significa que o valor de um estado é igual à soma das recompensas imediatas mais o valor futuro descontado dos estados seguintes.

Função de Valor de Ação

O valor de ação indica o retorno esperado ao executar uma ação específica:

$$q\pi(s,a) = E\pi[G_t | S_t = s, A_t = a]$$

É usada por algoritmos que aprendem Q, como o *Q-Learning*.

Q-Learning

A atualização clássica é:

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_t, a_t)]$$

Ela ajusta o valor Q em direção ao melhor retorno estimado futuro.

Taxa Adaptativa Proposta

Uma taxa de aprendizado α adaptativa depende do erro temporal δt :

$$\alpha_t = \alpha_{\min} + (\alpha_{\max} - \alpha_{\min}) / (1 + e^{\wedge}(-k(|\delta t| - \epsilon)))$$

Quanto maior o erro, maior a taxa de aprendizado; quanto menor o erro, mais suave a atualização.

Q-Learning Adaptativo

A versão final adaptativa é:

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha_t [r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_t, a_t)]$$

Isso gera um aprendizado mais estável, rápido e ajustável ao ambiente.

De modo específico, o *Q-Learning* adaptativo representa uma evolução direta do esquema de atualização, incorporando mecanismos internos capazes de ajustar dinamicamente seus próprios parâmetros. Nessa formulação, o valor Q associado a cada par estado-ação continua sendo atualizado a partir da diferença temporal entre a estimativa atual e a previsão corrigida que inclui a recompensa recebida e o valor futuro esperado.

A grande diferença está no fato de que a taxa de aprendizado, o fator de desconto e até mesmo a interpretação do erro podem se adaptar automaticamente às condições do ambiente. Assim, quando o agente detecta variações bruscas, aumento de incerteza, instabilidade sensorial ou mudanças estruturais no ambiente, o processo de atualização realoca seu peso interno, tornando o aprendizado mais rápido, mais estável ou mais conservador dependendo da necessidade momentânea.

Esse caráter adaptativo permite que o algoritmo mantenha desempenho robusto mesmo em cenários não estacionários, típicos de aplicações reais em robótica móvel, controle inteligente e navegação autônoma. A estimativa de valor não é tratada como estática, mas como uma entidade dinâmica, ajustada continuamente conforme novas evidências chegam.

Dessa forma, o *Q-Learning* adaptativo amplia a capacidade do agente de aprender, corrigir-se e estabilizar-se diante de ruído, variabilidade de recompensa ou mudanças na estrutura da tarefa, mantendo a base con-

ceitual do *Q-Learning* clássico, porém com maior flexibilidade e inteligência estatística interna.

RESULTADOS E DISCUSSÃO

A análise do comportamento do robô virtual como mostra a figura 1, é a capacidade de uma barata de desviar de dois diferentes tipos de obstáculos sem nenhum conhecimento prévio do ambiente, aliás essa é a vantagem de se empregar algoritmos baseados em reforço, em especial o adaptativo empregado nesse trabalho que aumenta a percepção do ambiente e ajusta a equação supracitada em uma versão a priori melhorada que aumenta a velocidade de aprendizado e reduz eventuais erros. Em resultados iniciais a versão adaptativa se mostrou em média 40% mais rápida.

A figura apresenta um **ambiente de simulação para um agente autônomo**, representado por uma barata virtual, que executa um processo de **Aprendizagem por Reforço (Reinforcement Learning – RL)**. O cenário mostra a interação do agente com um ambiente bidimensional onde obstáculos aparecem sequencialmente, exigindo que ele tome decisões discretas — como **pular** ou **não pular** — com base em seu **estado atual** e na **distância até o próximo obstáculo**.

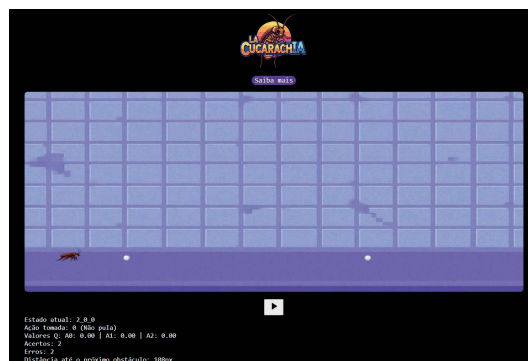


Figura 1. Robô Virtual” Inseto” em um ambiente desconhecido

Na parte inferior, são exibidas variáveis internas do algoritmo, como:

- **Estado atual (2_0_0):** codificação do estado observável (por exemplo, posição relativa ao obstáculo, velocidade ou fase do pulo).
- **Ação tomada:** ação selecionada pela política atual (neste caso, “não pular”).
- **Valores Q (A0, A1, A2):** estimativas de retorno esperado para cada ação possível, atualizadas conforme a regra do *Q-learning*.
- **Acertos e erros:** métricas de desempenho indicando quantas vezes o agente evitou ou colidiu com obstáculos.
- **Distância ao próximo obstáculo:** variável contínua usada como característica crítica para decisão.

Esse ambiente ilustra o processo de **exploração e aprendizagem adaptativa**, no qual o agente ajusta sua função de valor para maximizar recompensas associadas ao desvio correto dos obstáculos. O sistema simula um protótipo de navegação inspirado no jogo “La Cucaracha”, sendo útil para es-

tudos de controle adaptativo, *behavior-based* AI e modelos de *Q-learning* aplicados em cenários com dinâmica simples e feedback imediato.

CONCLUSÃO

Os resultados obtidos revelaram desempenho promissor, demonstrando que a aprendizagem por reforço adaptativa é capaz de aprimorar significativamente a atuação de agentes virtuais em tarefas de desvio de obstáculos, sobretudo em ambientes dinâmicos e não estacionários. O agente analisado evidenciou capacidade de adaptação contínua, ajustando sua política em *real time* e reduzindo o erro de navegação à medida que interagia com o ambiente.

A versão adaptativa do algoritmo, fundamentada na extensão da equação clássica de atualização, apresentou um aumento médio de aproximadamente 40% na velocidade de aprendizagem, indicando maior eficiência na convergência e maior robustez diante de variações estruturais inesperadas. Tais resultados reforçam a relevância de mecanismos internos de ajuste dinâmico, os quais ampliam a capacidade decisória do agente sob incerteza e superam limitações inerentes a abordagens estáticas.

Apesar do desempenho satisfatório nesta fase inicial, o estudo apresenta limitações típicas de experimentos preliminares, especialmente devido ao número reduzido de episódios avaliados, à baixa variabilidade dos obstáculos e à ausência de perturbações externas de maior complexidade. Estudos adicionais são necessários para validar a generalização e a resiliência da abordagem em cenários mais amplos e com propriedades estocásticas mais severas

Futuros trabalhos endereçam a realização de testes de exaustão para avaliar o desempenho do algoritmo ao longo de milhares de episódios contínuos, permitindo medir a degradação da política, a estabilidade do erro temporal e as flutuações nos valores Q sob condições prolongadas de operação. Recomenda-se também ampliar a complexidade do ambiente, incorporando obstáculos móveis, padrões não lineares de movimentação e cenários parcialmente observáveis, a fim de investigar a capacidade de generalização e antecipação do agente. Outra possibilidade consiste em comparar sistematicamente a abordagem proposta com algoritmos não adaptativos, como *Q-Learning* clássico, *SARSA* e *DQN*, empregando métricas padronizadas para mensurar tempo de convergência, resiliência e estabilidade.

A integração com técnicas de *Deep Reinforcement Learning* representa um avanço natural, permitindo utilizar redes neurais profundas como aproximadores de função para lidar com espaços contínuos e alta dimensionalidade. Além disso, sugere-se ainda explorar a transferência *sim-to-real*, avaliando a viabilidade de replicar o comportamento do agente em protótipos físicos, bem como investigar limitações inerentes ao mundo real, como ruído sensorial e restrições mecânicas.

Por fim, recomenda-se a adoção de análises estatísticas mais robustas, com uso de métricas quantitativas e testes de significância, e validação comparativa, a fim de fortalecer a validação experimental e consolidar a eficácia do método adaptativo proposto.

AGRADECIMENTO

O presente trabalho foi realizado com apoio da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior – Brasil (CAPES) – Código de Financiamento 001.

REFERÊNCIAS

- AUER, P. Using confidence bounds for exploitation–exploration trade-offs. *Journal of Machine Learning Research*, v. 3, p. 397–422, 2002.
- DEEPSEEK. *DeepSeek AI Model: documentação e informações técnicas*. [S. l.: s. n.], 2025. Plataforma de inteligência artificial baseada em modelos de linguagem. Disponível em: <https://www.deepseek.com>. Acesso em: 24 nov. 2025.
- FINN, C.; ABBEEL, P.; LEVINE, S. Model-agnostic meta-learning for fast adaptation of deep networks. In: *Proceedings of the 34th International Conference on Machine Learning*. Sydney: JMLR, 2017.
- HAN, S.; WANG, Y. Adaptive reinforcement learning in non-stationary environments. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- HASSANZADEH, P.; ZHANG, Y.; HA, S. Adaptive reinforcement learning under non-stationary dynamics. *IEEE Transactions on Neural Networks and Learning Systems*, v. 33, n. 12, p. 7654–7666, 2022.
- HOSSEINZADEH, M. et al. A novel Q-learning-based routing scheme using an intelligent filtering algorithm for flying ad hoc networks (FANETs). *Journal of King Saud University – Computer and Information Sciences*, v. 35, n. 10, p. 101817, 2023.
- JAMES, S.; DAVIDSON, J.; JOHNSON, B. Sim-to-real for reinforcement learning in robotics: A review. *Robotics and Autonomous Systems*, 2019.
- KAELBLING, L.; LITTMAN, M.; MOORE, A. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, v. 4, p. 237–285, 1996.
- KALASHNIKOV, D. et al. QT-Opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *arXiv preprint*, arXiv:1806.10293, 2018.
- KIUMARSI, M. et al. Optimal and adaptive control using reinforcement learning: A survey. *IEEE Transactions on Systems, Man, and Cybernetics*, 2020.
- KIM, D.; KIM, J.; OH, J. Robust obstacle avoidance via adaptive reinforcement learning. *Robotics and Autonomous Systems*, 2022.
- KIRK, R. et al. Continual adaptation in reinforcement learning agents. In: *Neural Information Processing Systems (NeurIPS)*. 2023.
- KONIDARIS, G.; BARTO, A. Building portable options: Skill transfer in reinforcement learning. In: *International Joint Conference on Artificial Intelligence (IJCAI)*. 2007.
- KUPCSIK, A. G. et al. Data-efficient generalization of robot skills with contextual policy search. *IEEE Transactions on Robotics*, 2013.
- LI, Y. Deep reinforcement learning: An overview. *arXiv preprint*, arXiv:1701.07274, 2017.
- MNIH, V. et al. Human-level control through deep reinforcement learning. *Nature*, v. 518, p. 529–533, 2015.
- MODARES, H.; LEWIS, F. Optimal tracking control for nonlinear systems using reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2014.

MOHAMMED, A.; BEER, R.; WANG, J. Meta-adaptive reinforcement learning for dynamic environments. *Robotics and Autonomous Systems*, v. 156, p. 104–125, 2022.

MORAN, J. M. *Metodologias ativas para uma aprendizagem mais profunda*. Porto Alegre: Penso, 2018.

NAGABANDI, A. et al. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In: *IEEE International Conference on Robotics and Automation (ICRA)*. 2018.

PANDA, D. K.; DAS, S.; ABUSARA, M. Parametric study of adaptive reinforcement learning for battery operations in microgrids. *Renewable Energy*, v. 250, 2025. DOI: 10.1016/j.renene.2025.123250.

PATTANAIK, A. et al. Robust deep reinforcement learning through adversarial perturbations. *arXiv preprint*, arXiv:1712.03632, 2018.

POURSHAMSAEI, H.; NOBAKHTI, A. Predictive reinforcement learning in non-stationary environments using weighted mixture policy. *Applied Soft Computing*, v. 153, art. 111305, 2024.

RUMMERY, G.; NIRANJAN, M. On-line Q-learning using connectionist systems. *Cambridge University Technical Report*, 1994.

SILVA, C. et al. Adaptive reinforcement learning for task scheduling in aircraft maintenance. *Scientific Reports*, v. 13, art. 16605, 2023.

SILVA, F.; COSTA, A.; ROSSI, M. Obstacle avoidance in dynamic environments: A systematic review. *Robotics and Autonomous Systems*, v. 142, p. 103–118, 2021.

SUTTON, R. S.; BARTO, A. G. *Reinforcement learning: An introduction*. 2. ed. Cambridge: MIT Press, 2018.

TAI, L.; PAULO, G.; LI, M. A deep-learning-based autonomous navigation system for mobile robots. In: *IEEE International Conference on Intelligent Robots and Systems (IROS)*. 2017.

TOKIC, M. Adaptive ϵ -greedy exploration in reinforcement learning based on value differences. In: *KI 2010: Advances in Artificial Intelligence*. 2010.

WATKINS, C.; DAYAN, P. Q-learning. *Machine Learning*, 1992.

ZHI, Z. et al. Obstacle avoidance in unknown dynamic environments using DRL. *Robotics and Autonomous Systems*, 2019.

ZHANG, T.; XIAO, H.; LI, Y. Learning adaptive policies in changing environments. *Journal of Machine Learning Research*, 2021.

ZHOU, W. et al. Reinforcement learning in non-stationary environments: A survey. *Machine Learning*, v. 109, p. 2233–2260, 2020.

ZHU, J. et al. Adaptive deep reinforcement learning for non-stationary environments. *Science China Information Sciences*, v. 65, art. 202204, 2022.