

# O que é Aprendizado por Reforço?

- Método de aprendizado de máquina onde um agente aprende a tomar decisões.
- Aprendizado por meio de tentativa e erro, interagindo com um ambiente.
- Recebe recompensas por ações corretas, visando maximizar a recompensa total.

# O que é Aprendizado por Reforço?



# Conceitos-chave

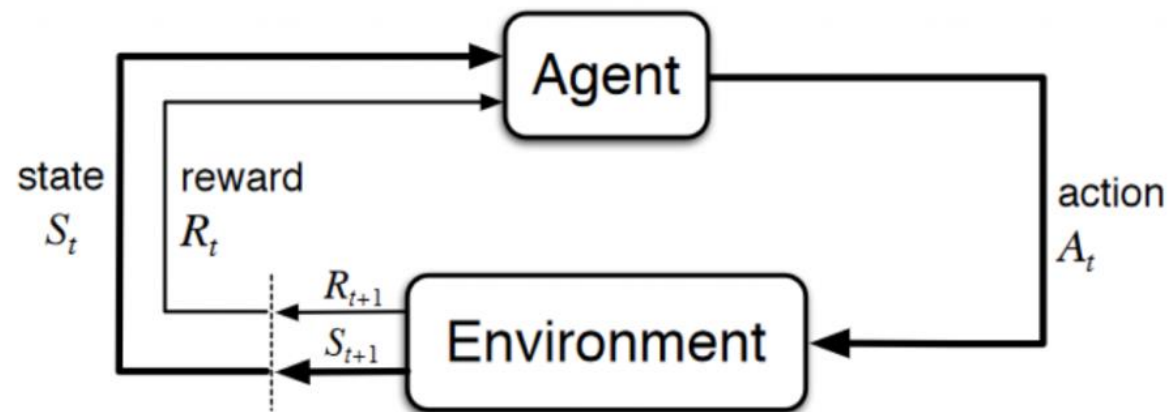
- **Agente:** A entidade que toma decisões, como um robô ou software.
- **Ambiente:** O mundo com o qual o agente interage e sobre o qual ele não tem controle total.
- **Estado:** Uma representação da situação atual do agente dentro do ambiente.
- **Ações:** As possíveis intervenções que o agente pode fazer no ambiente.

# Conceitos-chave

- **Recompensa:** Um sinal do ambiente em resposta às ações do agente, indicando o sucesso dessas ações.
- **Política (Policy):** Uma estratégia que define a escolha de ação do agente em determinados estados.
- **Função de Valor:** Uma previsão da recompensa futura esperada, utilizada para avaliar quão bom é um estado ou uma ação.
- **Episódio:** Uma sequência de ações, estados e recompensas que termina em um estado final.

# Como Funciona o Aprendizado por Reforço?

- Observação do estado do ambiente.
- Tomada de decisão baseada na política.
- Execução da ação e observação da recompensa e do novo estado.
- Atualização da política com base na recompensa recebida.



# Aplicações do Aprendizado por Reforço

- **Jogos:** Melhorar estratégias em jogos complexos, como Go e xadrez.
- **Robótica:** Ensinar robôs a realizar tarefas como caminhar e pegar objetos.
- **Sistemas de Recomendação:** Personalizar conteúdo para usuários em plataformas de streaming.
- **Otimização de Processos:** Melhorar a logística e a cadeia de suprimentos em indústrias.
- **Automação de Veículos:** Desenvolver sistemas de condução autônoma.

# Exploration x Exploitation

- **Exploração (Exploration)**
- Processo de tentar novas ações que o agente tem pouca ou nenhuma experiência anterior para descobrir informações valiosas sobre o ambiente.
- Essencial, especialmente em estágios iniciais do aprendizado, pois permite ao agente coletar dados sobre quais ações resultam em melhores ou piores recompensas.
- Sem exploração suficiente, um agente pode nunca descobrir estratégias ótimas ou mais eficazes para alcançar seu objetivo.
- A exploração vem com um custo, pois tentar ações desconhecidas pode levar a resultados negativos ou subótimos a curto prazo.

# Exploration x Exploitation

- **Exploração (Exploitation)**
- Ato de usar o conhecimento adquirido para tomar decisões que o agente já sabe que resultarão em boas recompensas.
- O desafio é que se um agente se concentrar exclusivamente na exploração, baseando-se apenas no conhecimento atual, ele pode perder a oportunidade de descobrir ações ainda mais recompensadoras que não foram suficientemente exploradas.



# Cenário de utilização

- Agente:
  - Um pequeno robô controlado por um programa de IA.
- Ambiente:
  - Um tabuleiro quadrado com várias células, algumas contendo obstáculos e uma contendo o tesouro.
- Recompensa:
  - O agente recebe uma recompensa positiva quando encontra o tesouro e uma recompensa negativa ao colidir com obstáculos ou sair do tabuleiro.
- Objetivo:
  - O objetivo do agente é aprender a política (estratégia) que maximiza sua recompensa cumulativa ao longo do tempo, ou seja, encontrar o tesouro com o mínimo de colisões possível.

# Cenário de utilização

- Inicialização:
  - O agente começa sem conhecimento sobre o ambiente.
  - Ele explora o ambiente fazendo ações aleatórias.
- Recompensas:
  - O agente recebe recompensas após cada ação.
  - Ele aprende que colidir com obstáculos ou sair do tabuleiro resulta em recompensas negativas, enquanto encontrar o tesouro resulta em uma recompensa positiva.
- Aprendizado da Política:
  - O agente usa um algoritmo de aprendizado por reforço para ajustar sua política com base nas recompensas recebidas.
  - Ele tenta maximizar as recompensas esperadas, aprendendo a evitar obstáculos e procurar o tesouro.

# Cenário de utilização

- Exploration vs. Exploitation:
  - O agente enfrenta o dilema de explorar novas ações (como tentar uma célula desconhecida) versus explorar ações conhecidas (como escolher uma célula que já sabe ser segura).
  - Ele deve encontrar o equilíbrio certo para maximizar a recompensa cumulativa.
- Aprimoramento Gradual:
  - Com o tempo, o agente aprimora sua política à medida que aprende a tomar decisões mais inteligentes com base em suas experiências passadas.

# Deep Q-Network (DQN)

- **Q-learning:**
- Visa aprender a política ótima, ensinando ao agente qual ação tomar sob determinadas condições para maximizar a soma de recompensas futuras.
- Isso é feito através de uma função de valor  $Q$ , que estima a recompensa total esperada de tomar uma ação em um dado estado.

# Deep Q-Network (DQN)

- **Redes Neurais Profundas:**
- No DQN, redes neurais profundas são usadas para aproximar a função de valor Q.
- Isso permite que o agente lide com estados de entrada de alta dimensão (como imagens de pixels dos jogos), algo que métodos tradicionais de Q-learning não conseguem fazer de forma eficiente devido à maldição da dimensionalidade.

# Questões de Concurso

Prova: CESPE / CEBRASPE - 2023 - SEFIN de Fortaleza - CE - Analista Fazendário Municipal - Área de Conhecimento: Ciência da Computação, Informática/Processamento de Dados

Julgue o item a seguir, a respeito de inteligência artificial (IA) e machine learning.

Nos algoritmos de aprendizado por reforço, o agente recebe uma recompensa atrasada na próxima etapa de tempo para avaliar sua ação anterior; seu objetivo, então, é maximizar a recompensa.

# Questões de Concurso

Prova: CESPE / CEBRASPE - 2023 - SEFIN de Fortaleza - CE - Analista Fazendário Municipal - Área de Conhecimento: Ciência da Computação, Informática/Processamento de Dados

Julgue o item a seguir, a respeito de inteligência artificial (IA) e machine learning.

Nos algoritmos de aprendizado por reforço, o agente recebe uma recompensa atrasada na próxima etapa de tempo para avaliar sua ação anterior; seu objetivo, então, é maximizar a recompensa.

# Questões de Concurso

Prova: CESPE / CEBRASPE - 2023 - DATAPREV - Analista de Tecnologia da Informação - Perfil: Inteligência da informação

Julgue o próximo item, relativos a aprendizado de máquina.

O aprendizado por reforço é um tipo de aprendizagem de máquina que tem por objetivo prever o resultado de um atributo alvo exclusivamente por meio de reforço no treinamento do modelo.



# Questões de Concurso

Prova: CESPE / CEBRASPE - 2023 - DATAPREV - Analista de Tecnologia da Informação - Perfil: Inteligência da informação

Julgue o próximo item, relativos a aprendizado de máquina.

O aprendizado por reforço é um tipo de aprendizagem de máquina que tem por objetivo prever o resultado de um atributo alvo exclusivamente por meio de reforço no treinamento do modelo.

# Questões de Concurso Inéditas

Qual das seguintes opções melhor descreve o princípio básico do aprendizado por reforço?

- A) O modelo é treinado exclusivamente com dados históricos, sem interação com o ambiente.
- B) O agente aprende a tomar decisões baseando-se unicamente em recompensas imediatas, sem considerar as consequências futuras.
- C) O agente aprende a tomar decisões através da experimentação no ambiente, buscando maximizar a soma de recompensas ao longo do tempo.
- D) O aprendizado ocorre por meio de instruções claras e diretas dadas ao agente antes da fase de testes.
- E) O agente utiliza um conjunto fixo de regras para tomar decisões, sem ajustes baseados na interação com o ambiente.

# Questões de Concurso Inéditas

Qual das seguintes opções melhor descreve o princípio básico do aprendizado por reforço?

- A) O modelo é treinado exclusivamente com dados históricos, sem interação com o ambiente.
- B) O agente aprende a tomar decisões baseando-se unicamente em recompensas imediatas, sem considerar as consequências futuras.
- C) O agente aprende a tomar decisões através da experimentação no ambiente, buscando maximizar a soma de recompensas ao longo do tempo.
- D) O aprendizado ocorre por meio de instruções claras e diretas dadas ao agente antes da fase de testes.
- E) O agente utiliza um conjunto fixo de regras para tomar decisões, sem ajustes baseados na interação com o ambiente.

# Questões de Concurso Inéditas

No contexto de aprendizado por reforço, o que é uma política (policy)?

- A) Uma estratégia fixa que o agente segue, determinada antes do início do treinamento.
- B) Um mapeamento de estados do ambiente para ações que o agente deve tomar, baseado em experiências passadas.
- C) A função que calcula a recompensa total acumulada pelo agente.
- D) Um registro de todas as ações tomadas pelo agente e as recompensas correspondentes.
- E) A configuração inicial do ambiente antes do agente começar a aprender.

# Questões de Concurso Inéditas

No contexto de aprendizado por reforço, o que é uma política (policy)?

- A) Uma estratégia fixa que o agente segue, determinada antes do início do treinamento.
- B) Um mapeamento de estados do ambiente para ações que o agente deve tomar, baseado em experiências passadas.
- C) A função que calcula a recompensa total acumulada pelo agente.
- D) Um registro de todas as ações tomadas pelo agente e as recompensas correspondentes.
- E) A configuração inicial do ambiente antes do agente começar a aprender.

# Questões de Concurso Inéditas

No aprendizado por reforço, o que é a função de valor?

- A) Uma previsão do número de ações necessárias para alcançar o objetivo.
- B) Uma estimativa das recompensas futuras que um agente pode esperar receber, estando em um determinado estado ou tomando uma certa ação.
- C) O total de recompensas que um agente recebeu até o momento.
- D) A probabilidade de um agente escolher a melhor ação em um dado estado.
- E) O custo computacional para realizar uma ação específica.

# Questões de Concurso Inéditas

No aprendizado por reforço, o que é a função de valor?

- A) Uma previsão do número de ações necessárias para alcançar o objetivo.
- B) Uma estimativa das recompensas futuras que um agente pode esperar receber, estando em um determinado estado ou tomando uma certa ação.
- C) O total de recompensas que um agente recebeu até o momento.
- D) A probabilidade de um agente escolher a melhor ação em um dado estado.
- E) O custo computacional para realizar uma ação específica.

# Questões de Concurso Inéditas

Qual dos seguintes algoritmos é um exemplo de aprendizado por reforço profundo (Deep Reinforcement Learning)?

- A) Linear Regression
- B) Decision Trees
- C) Deep Q-Network (DQN)
- D) K-Means Clustering
- E) Support Vector Machine (SVM)



# Questões de Concurso Inéditas

Qual dos seguintes algoritmos é um exemplo de aprendizado por reforço profundo (Deep Reinforcement Learning)?

A) Linear Regression

B) Decision Trees

C) Deep Q-Network (DQN)

D) K-Means Clustering

E) Support Vector Machine (SVM)

# Questões de Concurso Inéditas

Imagine um robô aprendendo (aprendizado por reforço) a caminhar em terrenos acidentados, onde cada passo dado é uma experimentação com o ambiente. O robô ajusta seus movimentos baseando-se em quais aspectos para melhorar seu desempenho ao longo do tempo?

- A) Seguindo um conjunto fixo de instruções programadas sem ajustes.
- B) Imitando movimentos humanos capturados previamente.
- C) Recebendo recompensas por manter o equilíbrio e penalidades por cair, ajustando seus passos de acordo.
- D) Analisando vídeos de outros robôs caminhando.
- E) Calculando a distância percorrida sem considerar a estabilidade.

# Questões de Concurso Inéditas

Imagine um robô aprendendo (aprendizado por reforço) a caminhar em terrenos acidentados, onde cada passo dado é uma experimentação com o ambiente. O robô ajusta seus movimentos baseando-se em quais aspectos para melhorar seu desempenho ao longo do tempo?

- A) Seguindo um conjunto fixo de instruções programadas sem ajustes.
- B) Imitando movimentos humanos capturados previamente.
- C) Recebendo recompensas por manter o equilíbrio e penalidades por cair, ajustando seus passos de acordo.
- D) Analisando vídeos de outros robôs caminhando.
- E) Calculando a distância percorrida sem considerar a estabilidade.

# Questões de Concurso Inéditas

Considere um sistema de IA treinado para jogar xadrez que começa sem conhecimento das estratégias do jogo. Ao longo de várias partidas, ele desenvolve a capacidade de prever movimentos do oponente e criar contraestratégias. Este processo exemplifica qual componente do aprendizado por reforço?

- A) Aprendizado supervisionado através de conjuntos de dados rotulados.
- B) Uso de uma base de dados de jogos históricos de xadrez para imitar grandes mestres.
- C) O agente melhora suas decisões ao maximizar suas recompensas, que, neste caso, são ganhar as partidas.
- D) Programação direta de estratégias de xadrez no agente.
- E) Aprendizado por imitação direta de movimentos de peças por humanos.

# Questões de Concurso Inéditas

Considere um sistema de IA treinado para jogar xadrez que começa sem conhecimento das estratégias do jogo. Ao longo de várias partidas, ele desenvolve a capacidade de prever movimentos do oponente e criar contraestratégias. Este processo exemplifica qual componente do aprendizado por reforço?

- A) Aprendizado supervisionado através de conjuntos de dados rotulados.
- B) Uso de uma base de dados de jogos históricos de xadrez para imitar grandes mestres.
- C) O agente melhora suas decisões ao maximizar suas recompensas, que, neste caso, são ganhar as partidas.
- D) Programação direta de estratégias de xadrez no agente.
- E) Aprendizado por imitação direta de movimentos de peças por humanos.

# Questões de Concurso Inéditas

Um agente de IA está aprendendo a otimizar rotas de entrega em uma cidade grande, onde o objetivo é minimizar o tempo total de entrega. A cada rota concluída, o agente recebe feedback que o ajuda a ajustar suas estratégias futuras. Que técnica de aprendizado por reforço é mais relevante para este cenário?

- A) Clustering para agrupar entregas por proximidade.
- B) Q-learning para ajustar a política de decisão com base no feedback recebido.
- C) Regressão linear para prever o tempo de entrega baseado na distância.
- D) Redes neurais convolucionais para processamento de imagens do mapa.
- E) Análise sentimental para entender o feedback dos clientes sobre as entregas.

# Questões de Concurso Inéditas

Um agente de IA está aprendendo a otimizar rotas de entrega em uma cidade grande, onde o objetivo é minimizar o tempo total de entrega. A cada rota concluída, o agente recebe feedback que o ajuda a ajustar suas estratégias futuras. Que técnica de aprendizado por reforço é mais relevante para este cenário?

- A) Clustering para agrupar entregas por proximidade.
- B) Q-learning para ajustar a política de decisão com base no feedback recebido.
- C) Regressão linear para prever o tempo de entrega baseado na distância.
- D) Redes neurais convolucionais para processamento de imagens do mapa.
- E) Análise sentimental para entender o feedback dos clientes sobre as entregas.

# Questões de Concurso Inéditas

Um software de IA, treinado para moderar comentários em um fórum online, aprende a identificar e remover postagens inapropriadas automaticamente. Inicialmente, comete erros, mas com o tempo, usando feedback dos usuários (recompensas e penalidades), ele melhora significativamente sua precisão. Este exemplo ilustra qual conceito-chave do aprendizado por reforço?

- A) A necessidade de um moderador humano para revisar todos os comentários.
- B) A importância de uma base de dados extensa de comentários previamente moderados.
- C) O uso exclusivo de técnicas de processamento de linguagem natural para entender o conteúdo dos comentários.
- D) A programação de regras específicas para cada tipo de comentário inapropriado.
- E) A habilidade do agente de aprender e otimizar seu comportamento através da interação com o ambiente.



# Questões de Concurso Inéditas

Um software de IA, treinado para moderar comentários em um fórum online, aprende a identificar e remover postagens inapropriadas automaticamente. Inicialmente, comete erros, mas com o tempo, usando feedback dos usuários (recompensas e penalidades), ele melhora significativamente sua precisão. Este exemplo ilustra qual conceito-chave do aprendizado por reforço?

- A) A necessidade de um moderador humano para revisar todos os comentários.
- B) A importância de uma base de dados extensa de comentários previamente moderados.
- C) O uso exclusivo de técnicas de processamento de linguagem natural para entender o conteúdo dos comentários.
- D) A programação de regras específicas para cada tipo de comentário inapropriado.
- E) A habilidade do agente de aprender e otimizar seu comportamento através da interação com o ambiente.