

TRABALHO DE GRADUAÇÃO

**APRENDIZAGEM POR REFORÇO NO CONTEXTO DE
FUTEBOL DE ROBÔS: *ROBOCUP SOCCER SIMULATION 2D***

Bruno Andregretti Dantas

Samuel Venzi Lima Monteiro de Oliveira

Brasília, Maio de 2021



**ENGENHARIA
MECATRÔNICA**
UNIVERSIDADE DE BRASÍLIA

UNIVERSIDADE DE BRASÍLIA
Faculdade de Tecnologia
Curso de Graduação em Engenharia de Controle e Automação

TRABALHO DE GRADUAÇÃO

**APRENDIZAGEM POR REFORÇO NO CONTEXTO DE
FUTEBOL DE ROBÔS: *ROBOCUP SOCCER SIMULATION 2D***

Bruno Andreghetti Dantas

Samuel Venzi Lima Monteiro de Oliveira

*Relatório submetido como requisito parcial de obtenção
de grau de Engenheiro de Controle e Automação*

Banca Examinadora

Prof. Dr. Alexandre Ricardo Soares Romariz
Orientador

Prof. Dr. Adolfo Bauschpiess
Examinador interno

Prof. Dr. João Paulo Leite
Examinador interno

Brasília, Maio de 2021

Dedicatórias

Dedico este trabalho aos meus pais, Chirlene e Roberto, e ao meu irmão, Pedro.

Samuel Venzi Lima Monteiro de Oliveira

Dedico este trabalho a todos aqueles que me trouxeram até aqui.

Bruno Andreghetti Dantas

Agradecimentos

Agradeço à minha família pela presença constante apesar da distância. Sem vocês minha jornada até aqui não seria possível. Agradeço também aos professores que me mostraram o caminho a ser trilhado e inspiraram a cada nova disciplina minha paixão pelo curso; aos meus colegas, que compartilharam comigo tantos dias e noites de estudo e também de jogatina por todos esses anos; aos amigos que fiz na DROID, pelas experiências que vivemos juntos no laboratório e nas competições; aos companheiros e ex-companheiros da GoLedger, onde evolui muito profissionalmente ao longo dos últimos anos; ao Samuel, que dividiu comigo essa jornada que completamos aqui; e, por fim, ao professor Romariz pela disponibilidade e atenção na orientação do trabalho.

Bruno Andregretti Dantas

Agradeço à minha família pelo papel fundamental na minha formação como indivíduo, por estar sempre presente para me divertir e me apoiar e pelo incentivo constante para que eu sempre buscasse fazer o melhor em todas as áreas da minha vida. Agradeço à Fernanda, minha namorada, pelas conversas e risadas e por compartilhar os melhores momentos comigo. Agradeço aos meus amigos que tornaram essa jornada mais leve: Victor Kzam, Ricardo Moura, André Luis, Felipe Caixeta, Gabriel Assis, Vieira Neto, Igor Beduin, Raphael Barbosa, Caio Campos, Vitor Duarte, Alice Lobo e Vitor Baltieri. Agradeço, também, à equipe UnBall por todos os desafios, noites em claro e viagens de competição. Agradeço ao Bruno, meu companheiro neste trabalho, pela amizade, competência e companheirismo nessa e em várias outras empreitadas. Finalmente, agradeço ao professor Alexandre Romariz pela orientação durante o desenvolvimento deste trabalho.

Samuel Venzi Lima Monteiro de Oliveira

RESUMO

O cenário da aprendizagem de máquina tem crescido cada vez mais nos últimos anos. Junto a isso, iniciativas como a RoboCup buscam incentivar a aplicação dessas técnicas fomentando um cenário competitivo de futebol de robôs. Com a motivação de aumentar a diversidade de ferramentas dentro da categoria *RoboCup Soccer Simulation 2D*, nesse trabalho foi desenvolvida uma nova plataforma de desenvolvimento e são aplicadas técnicas de aprendizagem por reforço a fim de validá-la. Foram realizados experimentos com técnicas estabelecidas como *Sarsa* e *Q-Learning* duplo tendo como objetivo realizar o maior número de gols possíveis durante o período de uma partida. As técnicas utilizadas validaram o funcionamento da plataforma desenvolvida e a utilização de comportamentos pré-programados aliado a *Q-Learning* duplo obteve uma política que alcançou gols consistentemente.

Palavras Chave: Aprendizagem por reforço, RoboCup, Futebol de robôs.

ABSTRACT

The machine learning field has been increasingly growing over the last few years. Furthermore, scientific initiatives such as RoboCup seek to promote studies and applications of these techniques by nurturing a competitive environment for robot soccer. In this project, with the motivation of expanding the diversity of tools in the RoboCup Soccer Simulation 2D category, a new development platform was created and reinforcement learning techniques were applied to validate it. Experiments were run using established techniques, such as Sarsa and Double Q-Learning, with the goal of scoring as many goals as possible during a match period. The techniques used validated the platform's operation and the use of pre-programmed behaviors along with Double Q-Learning resulted in a policy capable of consistently scoring goals.

Keywords: Reinforcement learning, RoboCup, Robot soccer

SUMÁRIO

1	Introdução.....	1
1.1	ROBÓTICA.....	1
1.2	APRENDIZAGEM POR REFORÇO.....	2
1.3	FUTEBOL DE ROBÔS.....	3
1.3.1	ROBOCUP SOCCER SIMULATION 2D.....	3
1.3.2	SERVIDOR DA PARTIDA.....	3
1.3.3	CLIENTE.....	4
1.3.4	SENSORES E AÇÕES.....	5
1.4	TRABALHOS RELACIONADOS.....	5
1.5	CARACTERIZAÇÃO DO PROBLEMA.....	6
1.5.1	OBJETIVOS ESPECÍFICOS.....	6
1.6	ESTRUTURA DO TRABALHO.....	7
2	Fundamentação Teórica.....	8
2.1	PROCESSOS DE DECISÃO DE MARKOV.....	8
2.1.1	MDP EPISÓDICO E CONTINUADO.....	9
2.1.2	RECOMPENSA E RETORNO.....	9
2.1.3	POLÍTICAS E FUNÇÕES DE VALOR.....	9
2.2	APRENDIZAGEM POR REFORÇO.....	10
2.2.1	APRENDIZAGEM COM DIFERENÇA TEMPORAL.....	10
2.2.2	APRENDIZAGEM <i>ON-POLICY</i> E <i>OFF-POLICY</i>	11
2.2.3	SOLUÇÕES TABULARES E APROXIMADAS.....	11
2.2.4	A TRÍADE DA MORTE.....	12
2.2.5	EXPLORAÇÃO E APROVEITAMENTO.....	12
2.2.6	SARSA.....	12
2.2.7	Q-LEARNING.....	13
2.2.8	Q-LEARNING DUPLO.....	13
3	Desenvolvimento.....	15
3.1	PLATAFORMA.....	15
3.1.1	ARQUITETURA DO CÓDIGO.....	15
3.1.2	DECODIFICAÇÃO E CODIFICAÇÃO DE MENSAGENS.....	15
3.2	DEFINIÇÃO DOS ESTADOS.....	16

3.2.1	SENSOR AUDITIVO	17
3.2.2	SENSOR VISUAL	17
3.2.3	SENSOR CORPORAL	18
3.2.4	ESTIMAÇÃO E TRATAMENTO DE ESTADOS	18
3.3	DEFINIÇÃO DE AÇÕES	19
3.3.1	ARRANCAR	20
3.3.2	CHUTAR	20
3.3.3	VIRAR.....	20
3.3.4	VIRAR PESCOÇO	20
3.3.5	MOVER-SE.....	20
3.3.6	FALAR.....	21
3.4	AMBIENTE DE TREINAMENTO	21
3.5	SARSA COM APROXIMADOR DE FUNÇÕES E COMPORTAMENTOS PRÉ-PROGRAMADOS	22
3.5.1	CODIFICAÇÃO DOS ESTADOS.....	22
3.5.2	CODIFICAÇÃO DOS COMPORTAMENTOS	23
3.5.3	PARÂMETROS	24
3.5.4	RECOMPENSA	24
3.6	<i>Q-LEARNING</i> DUPLO TABULAR E COMPORTAMENTOS PRÉ-PROGRAMADOS ...	25
3.6.1	CODIFICAÇÃO DOS ESTADOS.....	25
3.7	<i>Q-LEARNING</i> DUPLO TABULAR E AÇÕES ELEMENTARES	25
3.7.1	CODIFICAÇÃO DAS AÇÕES	26
3.7.2	PARÂMETROS	26
4	Resultados Experimentais.....	27
4.1	SARSA COM APROXIMADOR DE FUNÇÕES E COMPORTAMENTOS PRÉ-PROGRAMADOS	27
4.2	<i>Q-LEARNING</i> DUPLO TABULAR E COMPORTAMENTOS PRÉ-PROGRAMADOS ...	28
4.3	<i>Q-LEARNING</i> DUPLO TABULAR E AÇÕES ELEMENTARES	31
5	Conclusões.....	34
5.1	TRABALHOS FUTUROS.....	35
	REFERÊNCIAS BIBLIOGRÁFICAS	36

LISTA DE FIGURAS

1.1	Visualização de uma partida em andamento.....	4
1.2	Esquema ilustrando a arquitetura de um cliente e sua comunicação com o servidor do jogo.....	5
2.1	Interação agente-ambiente em um MDP [1].....	8
3.1	Indicadores espalhados pelo campo [2].	18
3.2	Visualização do cálculo da posição absoluta do jogador.	19
4.1	Curva de aprendizado do agente com comportamentos pré-programados utilizando Sarsa aproximado.	28
4.2	Curva de aprendizado do agente com comportamentos pré-programados.	29
4.3	Sobreposição de sequência de imagens do agente único fazendo gol.	29
4.4	Curva de aprendizado do agente com comportamentos pré-programados para treinamento longo.	30
4.5	Montagem ilustrando 4 exemplos da política aprendida através de sobreposição de imagens.....	31
4.6	Curva de aprendizado do agente com ações elementares.	32
4.7	Curva de aprendizado para treinamento longo.	32

LISTA DE TABELAS

3.1	Estrutura da rede neural utilizada no experimento.	22
-----	---	----

LISTA DE CÓDIGOS

3.1	Exemplo de mensagem codificada	16
3.2	Mensagem após passar pelo analisador léxico	16
3.3	Mensagem após passar pelo analisador sintático	16
3.4	Pseudo-código que descreve a rotina de treinamento	21
3.5	Algoritmo do comportamento de se aproximar da bola	23
3.6	Algoritmo geral do comportamento de carregar bola	23
3.7	Algoritmo geral do comportamento de chutar bola para o gol	24

LISTA DE SÍMBOLOS

Símbolos latinos

S	Estado
A	Ação
R	Recompensa
G	Retorno
v	Função de valor
q	Função de valor da ação
D	Medida de distância

Símbolos gregos

π	Política de tomada de ação
γ	Fator de desconto
α	Fator de aprendizagem
ϵ	Fator de exploração
θ	Medida de ângulo

Subscritos

t	Ciclo t de treinamento
$terminal$	Estado que encerra um episódio
π	Relativo a uma dada política π
$*$	Política ótima
j	Jogador
f	Bandeira do campo
$sensor$	Relativo ao sensor visual

Siglas

IA	Inteligência artificial
RL	Aprendizagem por reforço - <i>Reinforcement learning</i>
RCSS	<i>RoboCup Soccer Simulation 2D</i>
LARC	<i>Latin American Robotics Competition</i>
UDP	<i>User Datagram Protocol</i>
MDP	Processo de decisão de Markov - <i>Markov decision process</i>
GPI	<i>Generalized policy iteration</i>

Capítulo 1

Introdução

O treinamento de agentes, sejam eles máquinas físicas ou simuladas, para realizar certas tarefas é um campo de estudo crescente em meio à academia e ao setor privado. Em um ambiente de robôs jogadores de futebol simulados, este trabalho tem como alvo estruturar o procedimento de treinamento desses jogadores autônomos.

O interesse humano em criar artefatos para facilitar seu próprio trabalho ou realizar uma tarefa sem interferência remonta os tempos mais antigos. No Egito Ptolomaico, Ctesíbio (285-222 AC) descreveu um relógio d'água com a presença de um sistema de engrenagens, um indicador e o primeiro sistema de retroalimentação registrado. Por volta de 1495, Leonardo da Vinci, concebeu o projeto de um autômato mecânico de um guerreiro em armadura medieval que podia ficar em pé, sentar-se, levantar o visor e mover os braços [3].

O estudo da união de sistemas eletromecânicos e inteligência teve início há pelo menos 70 anos. A *cibernética*, área inaugurada por Norbert Wiener na década de 1950, descreve o estudo científico de controle e comunicação no animal e na máquina. Wiener começou, então, a desenvolver sistemas que replicassem comportamentos animais [4]. Somado a isso, a teoria da informação de Claude Shannon e a teoria de computação de Alan Turing abriram espaço para pesquisas que iriam desenvolver Inteligências Artificiais (IA) [5].

1.1 Robótica

A robótica se apoia em conhecimentos de vários campos para criar uma das áreas de estudos mais amplas da ciência. Desde a metade do século XX, a robótica vem reunindo noções dessas áreas e, pouco a pouco, tornando-as partes essenciais de si: sistemas mecânicos, eletromecânicos, teoria de controle, IA e outras. As aplicações existentes são inúmeras e se renovam a todo momento. Dentre as principais, é possível citar sistemas de manufatura, robótica médica e robótica agrícola [6].

Sistemas robóticos construídos com o objetivo de automatizar tarefas repetitivas se mostram úteis, entretanto o avanço das indústrias e aumento da complexidade das tarefas a serem realizadas

criou um ambiente catalisador para o desenvolvimento de processos capazes também de tomar decisões autonomamente no chão de fábrica.

É importante notar a complexidade do problema de se desenvolver a tomada de decisão de um sistema autônomo. Tal sistema precisa mapear seu ambiente por meio de sensores, extrair um significado do seu estado atual, usá-lo para decidir uma ação e determinar se tal ação foi a melhor a ser tomada. Sensores, porém, são imprecisos e limitados fisicamente. A representação dos estados, frequentemente, não é completamente conhecida e o processo de mapeamento de estado para ação não é trivial.

Neste contexto, surgiu a área de estudo conhecida como *aprendizagem por reforço*, que formaliza os elementos citados anteriormente para prover uma base de como agentes devem tomar ações para cumprir um objetivo pré-definido [1].

1.2 Aprendizagem por Reforço

A aprendizagem por reforço (do inglês, *reinforcement learning* ou RL) tem como inspiração a maneira como o aprendizado acontece com seres-humanos: interagindo com o ambiente [1]. Se uma criança está aprendendo a andar, por exemplo, ela toma certas ações no ambiente e, ainda que inconscientemente, está atenta aos resultados que essa ação causa.

A teoria por trás da aprendizagem por reforço formaliza a ideia de aprender através da interação e a aplica em um contexto computacional. Os principais elementos de um sistema de RL são: uma política de decisão, um sinal de recompensa e uma função valor. O primeiro diz respeito à decisão de qual ação se tomar a partir de uma situação, o segundo quantifica quão boa foi a ação escolhida naquele momento e o terceiro quantifica quão boa é a ação considerando o longo prazo [1].

Apesar de recente, a técnica de RL já se mostrou promissora em diversas áreas, com destaque para seu uso em jogos. Em 2016, o programa AlphaGo mostrou resultados significativos ao jogar contra o campeão europeu de Go, superando-o nos 5 jogos disputados [7].

Em 2018, pesquisadores do grupo OpenAI utilizaram técnicas de RL para treinar um time de 5 agentes colaborativos no jogo *Dota 2*, um jogo de estratégia em tempo real onde 2 times batalham para destruir a base inimiga. O jogo provê um ambiente extremamente complexo, com espaços de estados e ações contínuos. São 20.000 valores contínuos (ponto-flutuante), um para cada dimensão do estado, e 1.000 ações possíveis em um dado ciclo. A duração usual de uma partida é de pelo menos 1 hora. Inicialmente, foi feito um sistema do tipo um contra um (1v1) e o agente resultante deste treinamento foi capaz de derrotar jogadores profissionais. Em 2019, o sistema com 5 agentes foi capaz de derrotar um time profissional [8]. Um resultado dessa magnitude foi possível devido, entre outras razões, ao número altíssimo de amostras coletadas pelos agentes: 300 anos de experiência por dia para o agente singular e 180 anos por dia por agente para o time contendo 5 membros.

1.3 Futebol de Robôs

A ideia de robôs jogando futebol foi proposta pela primeira vez em 1992 por Alan Mackworth [9]. Desde então a comunidade científica tem criado iniciativas buscando por soluções que tornem isso realidade.

Uma delas é a *Robot World Cup Initiative* [10], abreviada como *RoboCup*, que teve sua primeira edição em 1997 com mais de 40 equipes distribuídas entre as diversas categorias do evento.

O objetivo da iniciativa, definido pela *RoboCup Federation*, é que por volta da metade do século XXI, um time de robôs humanóides autônomos vençam uma partida contra os campeões da Copa do Mundo mais recente. Mesmo que o objetivo pareça ambicioso, ele guia as pesquisas e motiva os avanços no campo.

Atualmente, a RoboCup conta com mais de 10 categorias, incluindo robôs humanóides, robôs com rodas e simulações. Entre elas há a *RoboCup Soccer Simulation 2D*, abreviada RCSS, objeto de estudo deste projeto.

1.3.1 RoboCup Soccer Simulation 2D

A RCSS possui grande relevância internacional, sendo uma das principais categorias disputadas na RoboCup, com equipes do mundo inteiro.

A categoria apresenta, também, grande relevância no cenário brasileiro. Desde 2005, a RCSS está presente na maior competição de robótica da América Latina, a *Latin American Robotics Competition*, LARC.

Nessa categoria, duas equipes de 11 jogadores autônomos e independentes jogam futebol em um ambiente virtual bidimensional. Um servidor é responsável por esse ambiente e possui informação absoluta sobre o estado do jogo e suas regras. Os jogadores, por sua vez, recebem dele informação incompleta e ruidosa de seus sensores virtuais, podendo executar comandos a fim de atuar sobre o estado do jogo [2].

1.3.2 Servidor da partida

Um servidor que executa a partida é disponibilizado pelos organizadores da competição e este pode ser utilizado, também, para desenvolvimento. O servidor, portanto, apresenta, internamente, algumas das regras da partida bem como um juiz autônomo que age para determinar gols, faltas e demais situações de uma partida de futebol. Caso necessário, um juiz humano poderá intervir em situações não contempladas pelas regras do servidor.

O servidor simula todos os movimentos e ações dos jogadores e da bola. Clientes externos se conectam ao servidor e cada cliente controla um único jogador. Uma partida usual tem 6.000 ciclos de 100 milissegundos, totalizando 10 minutos. A comunicação entre o cliente e o servidor é feita a partir do protocolo UDP por meio de mensagens com sintaxe específica e definida pelo servidor.

De forma a permitir o acompanhamento visual da partida, um monitor também é disponibilizado, porém não é necessário para que uma partida ocorra com sucesso.

O servidor possui, ainda, o modo *trainer* para utilização durante treinamentos de algoritmos de inteligência computacional. Este modo permite a conexão de um cliente do tipo treinador que tem acesso absoluto às informações da partida e pode mudar modos de jogo e ainda mover arbitrariamente jogadores e bola. Adicionalmente, é possível acelerar os ciclos da partida simulando o mais rápido possível e permitindo o treinamento em tempo hábil.



Figura 1.1: Visualização de uma partida em andamento.

1.3.3 Cliente

Os jogadores são controlados por clientes externos conectados ao servidor. Como citado, um cliente corresponde a um único jogador e os clientes só podem se comunicar entre si por meio de mensagens enviadas através do servidor da partida.

O cliente pode ser desenvolvido em qualquer linguagem desde que se comunique com o servidor via protocolo UDP e utilize a sintaxe de mensagens reconhecida pelo sistema. Há várias escolhas disponíveis para a construção do cliente, sendo decisão de cada equipe competidora como fazê-lo.

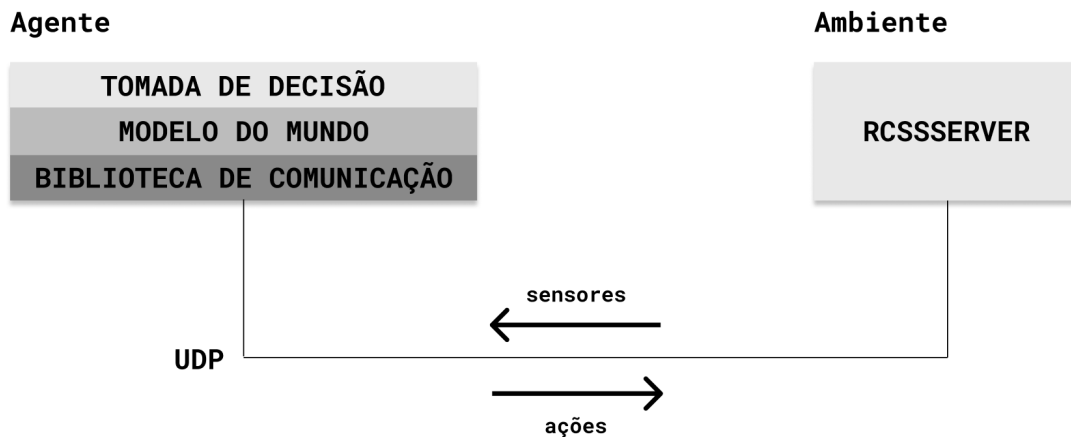


Figura 1.2: Esquema ilustrando a arquitetura de um cliente e sua comunicação com o servidor do jogo.

1.3.4 Sensores e Ações

Cada jogador presente na partida possui um conjunto de sensores de onde são tiradas todas as informações sobre o ambiente. Em uma partida usual, um jogador tem informações visuais dos jogadores do seu time e do time adversário, da bola e de uma série de marcadores fixos no campo, como bandeiras e linhas, que servem para situar o jogador em coordenadas absolutas do campo. O jogador possui também informações “sonoras”, onde pode ouvir mensagem do árbitro, treinador e de outros jogadores. Por último, têm acesso a informações do próprio corpo, como orientação do corpo e pescoço [2]. Os sensores possuem características que os aproximam de sensores reais como perda de resolução da informação conforme a variável medida se afasta do sensor.

Adicionalmente, a cada ciclo de simulação, cada cliente conectado ao servidor pode realizar ações que terão efeito no ambiente [2]. As ações englobam mover-se, virar-se, chutar a bola e até falar, permitindo troca de mensagens entre os jogadores. As ações disponíveis serão detalhadas no decorrer do texto.

1.4 Trabalhos Relacionados

Uma pesquisa sobre as abordagens para o desenvolvimento das estratégias dos times participantes da RCSS revelou o uso recorrente de métodos de inteligência computacional.

A equipe chinesa *WrightEagle*, campeã do principal evento internacional da categoria diversas vezes, utiliza Processos de Decisão de Markov ou MDPs para modelar a partida [11].

A equipe japonesa *HELIOS*, campeã de 2018 da categoria na RoboCup, divide seus jogadores em categorias “chutadores” e “não-chutadores”. Os chutadores são responsáveis por realizar o planejamento de sequência de ações, utilizando métodos de valor de ação. Os não-chutadores, por sua vez, não têm conhecimento do planejamento feito pelos chutadores e devem obter o máximo

de informações relevantes para tentar gerar a mesma sequência de ações que jogador chutador [12].

A equipe brasileira *ITAndroids*, atual campeã da LARC, utiliza a abordagem de sequência de ações, similar à *HELIOS*, explorando uma árvore de ações criada dinamicamente de forma a maximizar o valor de cada ação. Além disso, utilizam Otimização por Enxame de Partículas [13] para adequar os parâmetros que calculam o valor da ação. A *ITAndroids* também vem desenvolvendo o uso de Aprendizagem por Reforço Profunda [14].

Muitas equipes, ainda, desenvolvem seus agentes utilizando o agente base da equipe *HELIOS*, *Agent2d* com a biblioteca *Librcsc*, escritas em C++. O *Agent2d* utiliza uma abordagem que é chamada de *layered learning*, em que são definidas camadas nas quais são decididos comportamentos para cada agente dentro de estratégias pré-definidas [12]. Por ser utilizada por várias equipes, é comum que haja semelhança na construção entre agentes. Tal semelhança torna mais fácil a transferência de conhecimento entre as equipes, uma vez que todas utilizam a mesma plataforma. Porém, a reutilização das mesmas soluções acaba por limitar a diversidade de técnicas presentes nos campeonatos, o que pode prejudicar o avanço da área.

A utilização de comportamentos pré-programados, como os utilizados pela plataforma *Agent2d* é muito comum no âmbito do futebol de robôs, bem como em outras áreas onde o ambiente é conhecido e, portanto, é simples para um projetista codificar comportamentos a partir de seu conhecimento sobre o ambiente utilizando um conjunto de ações elementares. No ambiente de futebol de robôs, são exemplos de comportamentos: seguir bola e chutar bola para o gol. Nesse tipo de abordagem, o agente deve aprender a selecionar entre os comportamentos de acordo com a situação em que ele se encontra.

Entretanto, uma abordagem mais sofisticada seria utilizar as próprias ações elementares como objeto de escolha do agente para que haja mais flexibilidade no treinamento, permitindo o agente descobrir comportamentos não pensados pelo projetista. Exemplos de ações elementares seriam: andar, chutar e virar, onde parâmetros ditam como e com qual intensidade a ação será feita.

1.5 Caracterização do Problema

Deseja-se, então, explorar o problema de se fazer gols com um agente único no ambiente descrito utilizando-se de técnicas de aprendizagem por reforço. Para isso, foi desenvolvida uma biblioteca de interfaceamento com o servidor da partida como adaptação do ambiente. Após isso, o agente foi treinado com a utilização de técnicas de RL em duas abordagens - ações elementares e comportamentos - de forma a compará-las.

1.5.1 Objetivos Específicos

De acordo com o contexto apresentado, o presente trabalho se propõe a cumprir as seguintes etapas:

- Implementar uma plataforma para desenvolvimento de aprendizagem por reforço no âmbito da categoria RCSS;
- Utilizar técnicas de aprendizagem por reforço para treinar a escolha de comportamentos;
- Utilizar técnicas de aprendizagem por reforço para treinar a escolha de ações elementares;
- Comparar as diferentes configurações de treinamento e seus resultados.

1.6 Estrutura do Trabalho

Este documento está organizado em cinco capítulos: Introdução, Fundamentação Teórica, Desenvolvimento, Resultados e Conclusões. No Capítulo 2, descrevem-se os fundamentos teóricos da aprendizagem por reforço, com destaque para os algoritmos de *Sarsa*, *Q-learning* e *Q-learning* duplo. No Capítulo 3, descrevem-se o desenvolvimento da biblioteca e a definição de ações e comportamentos. Além disso, propõem-se o algoritmo de treinamento e os experimentos a serem realizados. No Capítulo 4, os resultados desses experimentos e suas análises são apresentados. Finalmente, o Capítulo 5 expõe conclusões e trabalhos futuros.

Capítulo 2

Fundamentação Teórica

2.1 Processos de Decisão de Markov

O problema abordado neste trabalho pode ser descrito simplificadaamente como um Processo de Decisão de Markov (MDP). MDP é uma forma clássica de representação matemática de processos de decisão sequenciais. Nesse modelo de representação, cada ação realizada por um agente que interage com o ambiente transforma o estado do processo e determina a recompensa que o agente recebe. Em um MDP, a codificação do estado atual deve conter toda a informação sobre a interação entre o agente e o ambiente que seja relevante para a dinâmica futura do processo. Nesse caso, é dito que o estado possui a "propriedade de Markov".

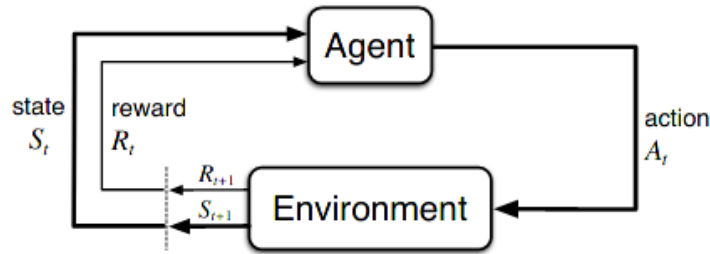


Figura 2.1: Interação agente-ambiente em um MDP [1].

Assim, dado um espaço de estados \mathcal{S} , um espaço de ações \mathcal{A} e um espaço de recompensas \mathcal{R} , para cada par (S, A) com $S \in \mathcal{S}$ sendo o estado atual do processo e $A \in \mathcal{A}$ a ação tomada pelo agente, existe uma determinada probabilidade de atingir o estado $S' \in \mathcal{S}$ e receber a recompensa $R \in \mathcal{R}$ [1].

Essa abordagem é bastante flexível e torna possível modelar a dinâmica do futebol virtual de robôs de diversas maneiras de modo que cada agente possa construir um estado percebido a partir de seus sensores e tomar decisões acerca de qual a melhor ação a fim de maximizar a recompensa recebida.

2.1.1 MDP Episódico e Continuado

Um MDP pode ser caracterizado quanto à presença de um estado terminal. Caso o MDP tenha um ou mais estados que determinem o fim do processo, ele é dito episódico e a totalidade do processo até o estado que determina o fim é chamada de episódio. A simulação de futebol de robôs tratada neste trabalho é um exemplo de MDP episódico, uma vez que o MDP termina ao se encerrar o tempo de jogo.

Em contrapartida, há MDPs nos quais não está bem definido nenhum estado terminal. Nesses casos, o MDP pode continuar indefinidamente até uma ação externa ao MDP determinar a sua parada. Um exemplo disso é um MDP que controle um robô numa linha de produção. Caso o sistema de automação supervisor desse robô não determine sua parada (por falta de insumos, por exemplo), o MDP pode seguir operando indefinidamente.

2.1.2 Recompensa e Retorno

Como definido acima, é atribuída uma recompensa $R \in \mathcal{R}$ para cada ação tomada em um MDP. Essa recompensa é sempre referente ao instante de tempo anterior, ou seja, não depende de qualquer outro fator que não o par (S_t, A_t) executado no instante t e a função de probabilidade associada pelo MDP a esse par. Por isso, é comum utilizar a notação R_{t+1} para se referir à recompensa obtida após tomar a ação A_t no instante de tempo t .

Porém, em muitos casos, é esperado de um agente que ele tome decisões para maximizar a recompensa total a longo prazo, ou seja, em um MDP episódico é esperado que se escolha A_t a fim de maximizar não apenas R_{t+1} mas sim o retorno G_t tal que:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_{terminal}. \quad (2.1)$$

Caso o MDP seja contínuo, G_t pode divergir, uma vez que é uma soma infinita de parcelas. Para solucionar o problema, basta adicionar à equação um fator de desconto (γ), tal que seja possível ajustar a relevância de recompensas futuras para o cálculo do retorno:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=1}^{\infty} \gamma^k R_{t+k+1}. \quad (2.2)$$

Observa-se que, para $0 \leq \gamma < 1$, o somatório converge. Para $\gamma = 0$, o retorno leva em consideração apenas a recompensa imediata R_{t+1} . Em contrapartida, para $\gamma \rightarrow 1$, as recompensas mais distantes no futuro são cada vez menos descontadas, ou seja, o agente tende a levar mais em consideração os ganhos futuros.

2.1.3 Políticas e Funções de Valor

É dado o nome de política para qualquer função $\pi : (S, A) \rightarrow \mathbb{R}$ que associe um estado qualquer do MDP e uma ação a uma probabilidade de se tomar essa ação diante desse estado, ou

seja, $\pi(A|S) = p$ onde p é a probabilidade de selecionar a ação A diante do estado S .

Para cada política π existe uma função $v_\pi : (S) \rightarrow \mathbb{R}$ que, para cada estado, define a esperança de retorno caso o agente continue seguindo a política π . A função v_π é conhecida como função de valor sob a política π .

Similarmente, existe uma função $q_\pi : (S, A) \rightarrow \mathbb{R}$ que, para cada par de estado e ação, define a esperança de retorno. Neste caso, a função q_π é conhecida como função de valor da ação sob a política π .

É possível comparar duas políticas π e π' a respeito de suas funções q . A política π é considerada melhor ou igual a π' , ou $\pi \geq \pi'$, caso $q_\pi(S, A) \geq q_{\pi'}(S, A)$ para todo par (S, A) .

Sempre há ao menos uma política melhor ou igual a todas as outras, denominada política ótima. Qualquer política que cumpra esse requisito é denominada π_* e, caso haja mais de uma, todas devem possuir a mesma função q denominada q_* [1]. No decorrer deste trabalho serão utilizadas técnicas que buscam estimar q_* .

Uma política que toma sempre o caminho de maior retorno é denominada gulosa e uma política que toma o caminho de maior retorno mas escolhe uma ação aleatoriamente com probabilidade parametrizada ϵ é denominada ϵ -gulosa.

2.2 Aprendizagem por Reforço

Dada uma modelagem do problema como um MDP, resta obter uma maneira de estimar as probabilidades que descrevem a dinâmica desse MDP para determinar a política capaz de maximizar a recompensa a longo prazo recebida pelo agente. O conjunto de técnicas que resolvem esse tipo de problema é chamado de Aprendizagem por Reforço.

No campo da aprendizagem de máquina, ela se difere da Aprendizagem Supervisionada por não haver um conjunto de pares (s, a) conhecidos a partir de exemplos. Nesse tipo de aprendizagem, o objetivo é extrapolar uma solução genérica a partir de exemplos de um conjunto de treinamento dado como correto, o que não é prático em problemas em que não se tem exemplos de comportamentos esperados e que representem bem o conjunto total de situações possíveis. Ela também se diferencia da Aprendizagem Não-Supervisionada, que tradicionalmente visa encontrar estrutura em conjuntos de dados não classificados, enquanto a Aprendizagem por Reforço visa maximizar um sinal de recompensa [1].

Desse modo, as técnicas de Aprendizagem por Reforço serão aplicadas a fim de buscar políticas capazes de maximizar o desempenho dos jogadores virtuais, ou seja, obter políticas que tornem os agentes capazes de fazer gols e evitar que os jogadores do time adversário façam gols.

2.2.1 Aprendizagem com Diferença Temporal

A aprendizagem com diferença temporal (*temporal-difference learning*, ou TD) é uma das ideias centrais de RL [1]. Esse tipo de aprendizagem é importante pois permite o aprendizado sem um

modelo prévio do MDP e também permite que se atualize as estimativas das funções de valor a cada ciclo de treinamento utilizando a estimativa anterior como aproximação, técnica conhecida como *bootstrapping*, que reduz o tempo de treinamento e viabiliza a aprendizagem em MDPs continuados.

A diferença temporal se apoia na ideia de que a política aproximada e a função de valor devem se relacionar de forma que ambas convirjam para seus valores ótimos, conhecida como iteração generalizada da política (*generalized policy iteration* ou GPI) [1].

As técnicas *Sarsa*, *Q-learning* e *Double Q-learning* apresentadas respectivamente nas Subseções 2.2.6, 2.2.7 e 2.2.8 são exemplos de aprendizagem com diferença temporal.

2.2.2 Aprendizagem *On-policy* e *Off-policy*

Entre as técnicas de aprendizagem por reforço, existe uma divisão entre a aprendizagem *on-policy* e a aprendizagem *off-policy*, referentes à relação entre a política executada durante o aprendizado e a política sobre a qual se quer aprender.

Nos algoritmos de aprendizagem *on-policy* o agente aprende a respeito da política π enquanto navega o MDP de acordo com a própria política π , ou seja, a política executada durante a aprendizagem é a mesma que se quer estudar.

Em contrapartida, nos algoritmos de aprendizagem *off-policy* o agente aprende a respeito da política alvo π enquanto navega o MDP de acordo com a política b , ou seja, ele estima a função q_π enquanto executa a política b . Esses métodos costumam introduzir variância no processo, tornando o aprendizado ruidoso e muitas vezes divergente.

Além disso, é possível observar que a aprendizagem *on-policy* é apenas um caso particular da aprendizagem *off-policy* em que $b = \pi$.

2.2.3 Soluções Tabulares e Aproximadas

A maioria dos métodos de aprendizagem por reforço são testados e validados em MDPs cujos espaços de estados \mathcal{S} e de ações \mathcal{A} são suficientemente pequenos. Para esses MDPs é possível utilizar uma solução tabular, ou seja, a função Q pode ser armazenada em uma tabela de tamanho razoável e sua imagem para cada par estado-ação pode ser atualizado individualmente.

Infelizmente, em diversas aplicações a quantidade de estados possíveis é grande demais ou até mesmo infinita, como é o caso de sistemas em que determinada característica do estado é medida como uma grandeza contínua. Nesses casos, é impossível esperar que se obtenham soluções ótimas mesmo com tempo infinito, portanto o objetivo é obter uma solução aproximada que seja boa o suficiente para a aplicação desejada.

A ferramenta matemática utilizada para viabilizar soluções aproximadas é o conceito de aproximadores de função, muito utilizados na aprendizagem supervisionada. Entre os aproximadores mais utilizados, estão os aproximadores lineares e as redes neurais multicamada.

2.2.4 A Tríade da Morte

Sutton e Barto definem um cenário denominado a tríade da morte (*the deadly triad*) como um cenário onde a presença de três elementos causa instabilidade e divergência em treinamentos: aproximação de funções, *bootstrapping* e treinamentos *off-policy* [1].

Existe, portanto, a decisão do projetista de definir treinamentos de forma a evitar o conjunto desses três elementos. O aproximador de funções pode ser evitado com técnicas tabulares em ambientes menos complexos e ao custo de uso de memória. Para métodos TD, *bootstrapping* é um requerimento. Entretanto, outros métodos permitem realizar o treinamento sem o uso de estimativas. Métodos *off-policy* são considerados mais poderosos e versáteis, porém introduzem mais variância como descrito na Subseção 2.2.2 e, portanto, nem sempre são a melhor escolha.

2.2.5 Exploração e Aproveitamento

Na realização do treinamento de um agente por RL, surge o dilema da exploração e do aproveitamento (na literatura em inglês chamados de *exploration* e *exploitation* respectivamente). Para que o agente descubra as ações que geram os maiores retornos, é necessária a ampla experimentação de pares (S, A) disponíveis no MDP. Sem experimentar novas ações e visitar novos estados, não é possível que o agente descubra o melhor caminho.

Oposto a isso, para obter um maior retorno, o agente deve dar preferência às ações já visitadas nas quais se sabe que produz recompensas maiores, aproveitando o conhecimento adquirido anteriormente.

É necessário, portanto, balancear a exploração e o aproveitamento a fim de direcionar o treinamento. Uma das técnicas utilizadas para isso é a adoção de uma política ϵ -gulosa, descrita na Seção 2.1.3, que permite ajustar o equilíbrio entre exploração e aproveitamento através do parâmetro ϵ .

Além disso, alguns ambientes permitem que se escolha o estado inicial do agente. Nesse caso, desde que o ambiente seja episódico, é possível selecionar aleatoriamente o estado inicial a cada novo episódio e dessa forma garantir que o agente visite novos estados independentemente da política adotada.

2.2.6 Sarsa

Entre os algoritmos de aprendizagem *on-policy* está o Sarsa. Seu nome é derivado da quintupla $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$ utilizada como entrada da sua fórmula de atualização [1]:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]. \quad (2.3)$$

Uma vez que a política adotada pelo agente não pode ser gulosa, pois inibiria a exploração do espaço de estados e de ações, é interessante que o agente siga uma política ϵ -gulosa baseada na

estimativa de Q . Nesse caso, para Q tabular, a convergência para a política ótima é garantida desde que todos os pares (S, A) sejam visitados infinitas vezes e a política convirja para a política gulosa em $t \rightarrow \infty$.

Para o caso em que se usa um aproximador para estimar Q não há garantias, porém a natureza *on-policy* do Sarsa reduz a variância da aprendizagem. Isso pode tornar essa abordagem mais estável do que o Q-learning, descrito a seguir.

2.2.7 Q-Learning

Um dos algoritmos mais populares no campo da aprendizagem por reforço é o Q-Learning. Trata-se de um método *off-policy* que aproxima diretamente a função q_* independente da política que estiver sendo adotada pelo agente durante o treinamento.

O algoritmo é também muito simples. Dada uma representação tabular $Q : (\mathcal{S}, \mathcal{A}) \rightarrow \mathbb{R}$ da função q_* , para cada instante de tempo t é realizada a seguinte atualização a fim de aproximar Q de q_* [1]:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)], \quad (2.4)$$

sendo α o fator de aprendizagem, responsável por suavizar o impacto de cada atualização da tabela, e γ o fator de desconto, responsável por reduzir a relevância de recompensas muito distantes no tempo.

Após iterações suficientes, espera-se que Q convirja para q_* . Em certas condições a convergência é matematicamente garantida.

Uma vez estimada a função q_* , é simples obter a política ótima. Basta escolher a ação que maximiza q_* no estado atual, ou seja:

$$A_t = \max_a q_*(S_t, a). \quad (2.5)$$

É comum, mas não obrigatório, que a política b seguida durante o aprendizado seja *epsilon*-gulosa em relação à aproximação Q .

2.2.8 Q-Learning Duplo

Apesar de popular, o Q-Learning possui um problema de viés de maximização. Uma vez que a aproximação Q é imprecisa no início do treinamento, é possível que o retorno esperado estimado seja enviesado para um valor maior do que o real.

Como solução para esse problema, é utilizada a abordagem do Q-Learning Duplo. Nela são utilizadas duas aproximações, Q_1 e Q_2 , e a atualização de Q é dada da seguinte forma [1]:

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha[R_{t+1} + Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t)]. \quad (2.6)$$

Em metade das iterações (através de um sorteio, por exemplo), as aproximações Q_1 e Q_2 são trocadas. Com isso, é anulado o viés de maximização gerado pelo uso de $\max_a Q$ como estimativa de retorno para os estados seguintes.

Outra vantagem desse método é que, apesar de dobrar os requisitos de memória do algoritmo por ser preciso armazenar os dados referentes a duas aproximações, ele não aumenta o custo computacional por iteração.

No capítulo seguinte estão descritos o desenvolvimento da plataforma de testes e também os experimentos realizados para validá-la.

Capítulo 3

Desenvolvimento

3.1 Plataforma

O servidor da partida apresenta, como já mencionado na Subseção 1.3.2, um protocolo de comunicação e sintaxe de mensagens específicos. Uma plataforma de interfaceamento foi desenvolvida com o objetivo de abstrair os detalhes de comunicação e de construção de mensagens e facilitar, assim, o desenvolvimento dos programas jogadores. Esta abordagem já é comum na categoria e existem soluções de código aberto como a *librcsc*, utilizada por várias equipes, usualmente atreladas ao agente base *agent2d*, desenvolvidas pela equipe *HELIOS*, como visto na Subseção 1.4.

Uma biblioteca própria foi desenvolvida em linguagem Go como forma de modernização e diversificação da base de código utilizada pelas equipes. A biblioteca cobre uma parte considerável das possibilidades previstas no protocolo de comunicação e foi programada de modo a ser facilmente expansível de acordo com o lançamento de atualizações do servidor.

3.1.1 Arquitetura do código

A biblioteca possui três pacotes internos: *playerclient*, *trainerclient* e *rcsscommon*. Os dois primeiros dizem respeito aos dois tipos de programas que podem se conectar ao servidor da partida: jogadores e treinadores. O terceiro engloba todas as funcionalidades utilizadas por ambos clientes, além de informações gerais sobre parâmetros da partida, como coordenadas de bandeiras do campo e modos de jogo. Os dois clientes desenvolvidos possuem as funcionalidades necessárias para se conectar ao servidor, ouvir mensagens via protocolo UDP, decodificá-las e então executar uma ação em forma de mensagem codificada e enviada ao servidor.

3.1.2 Decodificação e Codificação de Mensagens

A decodificação de mensagens, por sua vez, foi feita em duas camadas: um analisador léxico e um analisador sintático. O analisador léxico passa pelas mensagens em formato *string* e identifica os símbolos que ela contém. O analisador sintático extrai desses símbolos as informações e as

organiza em estruturas de dados para que possam ser utilizadas fora da biblioteca. As informações recebidas e decodificadas são, em sua maioria, dados dos sensores do jogador.

```
(see 37 ((f c b) 16.6 -1 -0 -0.8))
```

Código 3.1: Exemplo de mensagem codificada

```
SightSymbols{
  Time: 37,
  ObjMap: map[string][]string{
    "f c b": {"16.6", "-1", "-0", "-0.8"},
  },
}
```

Código 3.2: Mensagem após passar pelo analisador léxico

```
SightData{
  Time: 37,
  Ball: nil,
  Lines: LineArray{},
  Flags: FlagArray{
    {
      ID: rcsscommon.FlagCenterBot,
      Distance: 16.6,
      Direction: -1,
    },
  },
}
```

Código 3.3: Mensagem após passar pelo analisador sintático

3.2 Definição dos Estados

As informações de estado são fornecidas pelos sensores do jogador. Há três sensores presentes que entregam informações em forma de mensagens para o agente: sensor auditivo, sensor visual e sensor corporal.

Além de ser necessário decodificar as mensagens recebidas, como demonstrado na Subseção 3.1.2, é necessário processar os dados dos sensores para extrair mais características do estado atual do agente.

Os sensores são responsáveis por todas as informações que o jogador tem do ambiente. Eles são modelados de forma a emular características de sensores reais, portanto um sensor pode “perder” informações caso o objeto da variável medida esteja longe e, além disso, ruído é injetado ao valor medido dessa variável, como será evidenciado no modelo do sensor visual.

3.2.1 Sensor Auditivo

As mensagens do sensor auditivo são do seguinte formato:

(hear Tempo Remetente "mensagem"),

onde *Tempo* é o número do ciclo em que a mensagem foi ouvida e *Remetente* é descrição de quem enviou a mensagem. O *Remetente* pode ser o árbitro, outros jogadores, um dos treinadores ou o próprio jogador.

No escopo deste trabalho, apenas as mensagens do árbitro serão consideradas, não sendo implementada nenhuma forma de comunicação direta entre os agentes.

3.2.2 Sensor Visual

As mensagens do sensor visual contém as posições relativas referentes a cada objeto dentro do campo de visão do jogador. Esses objetos podem ser outros jogadores, marcadores como bandeiras e linhas (Figura 3.1) e a bola. O formato genérico é este:

(see (Objeto1)(Objeto2)(Objeto3)...(ObjetoN)),

onde cada objeto tem o seguinte formato:

((NomeDoObjeto) Distância Direção VariaçãoDeDistância VariaçãoDeDireção DireçãoDoCorpo DireçãoDaCabeça),

sendo a distância e a direção dados em coordenadas polares, assim como suas variações. As informações de direção do corpo e da cabeça só aparecem quando o objeto em questão é um outro jogador.

O campo é marcado com vários indicadores com posição fixa conhecida para que o agente possa estimar sua posição absoluta, como vista na Figura 3.1.

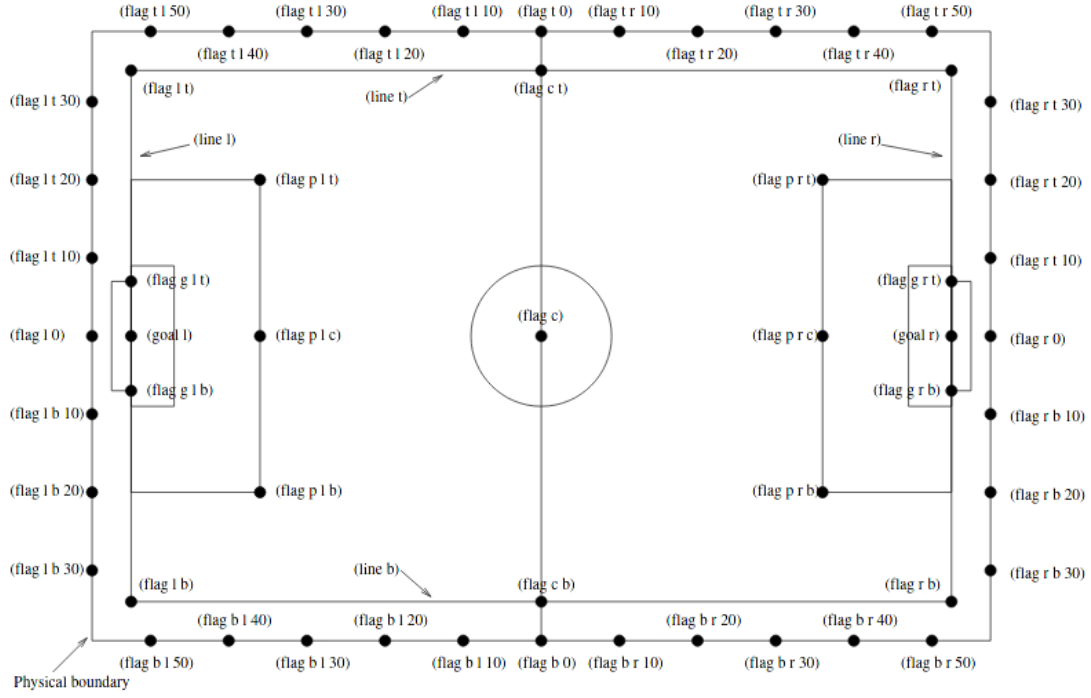


Figura 3.1: Indicadores espalhados pelo campo [2].

A riqueza de detalhes a respeito das informações obtidas depende da distância entre o objeto e o jogador. Por exemplo, caso um outro jogador esteja sendo visto a uma distância muito grande, pode ser impossível determinar o número de sua camisa ou até mesmo a qual time ele pertence. Em contrapartida, para jogadores próximos, é fornecido até mesmo a direção para a qual ele está olhando.

3.2.3 Sensor Corporal

O sensor corporal contém informações sobre o estado físico do jogador. Entre elas sua energia, que é consumida a cada ação tomada como chute ou arrancada (Subseção 1.3.4), sua própria velocidade e direção de movimento, a direção de sua cabeça e a quantidade de cartões de advertência recebidos.

3.2.4 Estimação e Tratamento de Estados

Os sensores do jogador fornecem somente informações em coordenadas polares relativas ao próprio jogador. Desta forma, ele possui informações de distância e direção para a bola, demais jogadores, bandeiras do campo e linhas do campo.

Com esses dados é possível estimar as coordenadas cartesianas absolutas no campo de todas as entidades de interesse. As bandeiras do campo são fixas e possuem coordenadas conhecidas, como detalhado na subseção 3.2.2. Logo, a transformação da informação polar e relativa para uma cartesiana e absoluta da posição do próprio jogador pode ser realizada.

Inicialmente, obtém-se o ângulo θ_j do jogador em relação ao eixo X a partir do ângulo das linhas no perímetro do campo. Basta subtrair o ângulo absoluto da linha em relação ao ângulo medido pelo sensor.

De posse do ângulo absoluto do jogador, sua posição absoluta (X, Y) em campo é calculada da seguinte forma a partir das bandeiras vistas:

$$\begin{aligned}\theta_{j \rightarrow f} &= \theta_{sensor} + \theta_j; \\ X &= X_f - D_{sensor} \cos(\theta_{j \rightarrow f}); \\ Y &= Y_f - D_{sensor} \sin(\theta_{j \rightarrow f}),\end{aligned}\tag{3.1}$$

sendo D_{sensor} e θ_{sensor} a distância e ângulo medidos pelo sensor, $\theta_{j \rightarrow f}$ o vetor que leva do jogador à bandeira e (X_f, Y_f) a posição absoluta da bandeira. Uma média é feita a partir das estimativas obtidas com cada uma das bandeiras vistas.

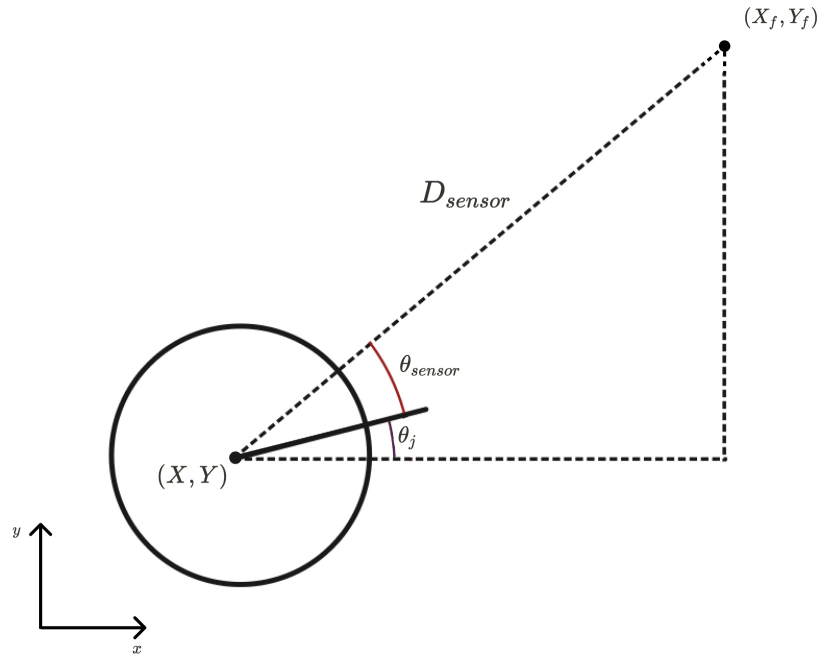


Figura 3.2: Visualização do cálculo da posição absoluta do jogador.

A partir da informação de posição absoluta do próprio jogador, é possível calcular as posições absolutas para o resto das entidades observadas.

3.3 Definição de Ações

Os jogadores possuem diversas ações possíveis mapeadas pelo servidor da partida. As ações recebem parâmetros em sua maioria contínuos. Dessa forma a escolha de quais ações devem ser executadas se dá em um domínio discreto, porém cada ação exige uma escolha de parâmetros em um domínio contínuo.

3.3.1 Arrancar

O comando de arrancar faz com que haja uma aceleração do jogador na direção da arrancada. O parâmetro “potência” determina o valor da aceleração e o parâmetro “direção” é relativo à direção do corpo do jogador. É importante salientar que o comando de arranque é a forma padrão de movimentar um jogador.

Cada jogador possui uma certa quantidade de energia e o arranque tem um custo sobre ela. Ao começo de cada partida, a energia do jogador é colocada no máximo. Se o jogador acelera para frente, a energia é reduzida em $1 \times$ a potência. Se o jogador acelera para trás, o custo é maior e a energia é reduzida em $2 \times$ a potência.

Se a energia disponível é menor que a necessária para a realização com comando, o valor de “potência” é reduzido para que a quantidade necessária de energia seja a disponível.

3.3.2 Chutar

O comando de chute recebe dois parâmetros: a força e a direção do chute. Para realizar o comando, a bola precisa ser “chutável”, ou seja, estar a uma certa distância do jogador.

Caso a bola não esteja diretamente à frente do jogador, a força efetiva será reduzida por um fator dependente da posição relativa da bola.

3.3.3 Virar

O comando virar recebe como parâmetro o momento angular a ser aplicado pelo jogador sobre si mesmo. Se o jogador não estiver em movimento, o seu ângulo é apenas incrementado com o momento.

Porém, caso o jogador esteja em movimento, o resultado do comando é influenciado pelo momento de inércia do jogador (definido aleatoriamente pelo servidor no início da partida) e sua velocidade linear.

3.3.4 Virar pescoço

O jogador pode virar seu pescoço de maneira semi-independente de seu corpo. O ângulo do pescoço é relativo ao ângulo do corpo, então caso um comando virar seja executado, o ângulo absoluto do pescoço também mudará. Os ângulos máximo e mínimo em relação ao corpo são definidos na configuração do servidor, sendo o padrão $+90$ e -90 graus, respectivamente.

3.3.5 Mover-se

O comando “mover-se” é utilizado para mover diretamente jogadores para coordenadas. Não pode ser utilizado no decorrer de uma partida, exceto quando o goleiro tem a posse da bola. O

comando fica disponível no início dos tempos da partida para posicionamento dos jogadores. Para movimentar os jogadores durante a partida, o arranque (Subseção 3.3.1) deve ser utilizado.

3.3.6 Falar

Usado para transmitir mensagens aos outros jogadores. As mensagens devem ter comprimento menor que um valor pré-determinado pelo servidor. Os jogadores que estiverem a uma distância audível da mensagem receberão a mensagem do servidor imediatamente.

Não será implementada nenhuma comunicação entre os jogadores neste trabalho.

3.4 Ambiente de Treinamento

O ambiente de treinamento consiste em uma base de código que importa a biblioteca detalhada na Seção 3.1. Um formato geral foi definido e desenvolvido a fim de tornar os experimentos fáceis de adaptar, bastando mudar alguns trechos do código. No contexto desse trabalho, um episódio representa 6.000 ciclos com um único jogador e a bola em campo.

O laço de treinamento obedece ao pseudo-código apresentado no Código 3.4.

```
definir parâmetros
inicializar pesos de treinamento
enquanto o estado não é terminal:
  conectar jogador
  para cada ciclo do episódio:
    escolher ação de acordo com a política e o estado
    observar novo estado e recompensa
    treinar pesos
  estado <- novo estado
```

Código 3.4: Pseudo-código que descreve a rotina de treinamento

O laço interno é onde efetivamente os algoritmos de treinamento são implementados, portanto este trecho é alterado a depender da técnica de aprendizagem por reforço utilizada.

O objetivo inicial para teste do sistema como um todo foi o de realizar o treinamento de um agente capaz de executar gols estando sozinho em campo.

Para isso, foram realizadas três abordagens diferentes.

- *Sarsa* com aproximador de funções e comportamentos pré-programados;
- *Q-Learning* duplo tabular e comportamentos pré-programados;
- *Q-Learning* duplo tabular e ações elementares.

3.5 Sarsa Com Aproximador de Funções e Comportamentos Pré-Programados

Inicialmente, desejou-se realizar um treinamento utilizando um método aproximado, devido ao grande número de estados contínuos possíveis para o agente. Portanto, foi usada uma rede neural multicamada com estrutura detalhada na Tabela 3.1.

A função de ativação utilizada foi *ReLU*. A camada de entrada da rede apresenta 5 nós, cada um relativo a um dos estados descritos na Subseção 3.5.1, e a camada de saída apresenta 8 nós, cada um relativo a um dos comportamentos descritos na Subseção 3.5.2. A cada ciclo t , o agente toma uma ação A_t no estado S_t e recebe o próximo estado S_{t+1} e recompensa R_{t+1} . A partir disso, o valor $Q(S_t, A_t)$ é calculado de acordo com a equação 2.3 e, então, o treinamento dos pesos da rede é realizado de forma a aproximar o valor calculado.

camada	nº nós
1	5
2	256
3	256
4	128
5	64
6	8

Tabela 3.1: Estrutura da rede neural utilizada no experimento.

A utilização de comportamentos permite que sejam inseridos conhecimentos prévios a respeito do que é esperado de um agente jogador de futebol, ou seja, é possível simplificar o aprendizado substituindo as ações elementares, como chutar ou correr, por ações de mais alto nível de abstração, os comportamentos, como perseguir a bola e chutar para o gol.

3.5.1 Codificação dos Estados

- **Coordenada X do Jogador:** dentro do intervalo $[-1, 1]$ e normalizada em relação ao comprimento do campo.
- **Coordenada Y do Jogador:** dentro do intervalo $[-1, 1]$ e normalizada em relação à largura do campo.
- **Orientação do Jogador:** dentro do intervalo $[-1, 1]$ e proporcional ao ângulo do corpo do jogador em relação ao eixo X entre -180° e 180° .
- **Distância até a bola:** normalizada em relação à distância máxima visível.
- **Direção da bola:** dentro do intervalo $[-1, 1]$ e proporcional ao ângulo do corpo do jogador em relação ao vetor que o liga à bola entre -180° e 180° .

3.5.2 Codificação dos Comportamentos

Os comportamentos pré-programados foram desenvolvidos de forma a cobrir ações que permitissem a realização de gols, que é o objetivo do treinamento.

- **Comportamento nulo:** O agente apenas espera até o próximo ciclo.
- **Localizar bola:** Caso o agente não veja a bola, ele gira no seu próprio eixo na direção em que a bola foi vista pela última vez. Se a bola não foi vista por mais de 30 ciclos, o agente gira em seu eixo 45° , a fim de achar a bola mais rápido.
- **Se aproximar da bola:** Caso o agente veja a bola, ele se move em direção a ela com velocidade constante
- **Se afastar da bola:** Caso o agente veja a bola, ele se move em direção oposta a ela com velocidade constante
- **Carregar bola para o gol direito:** Caso a distância até a bola seja menor ou igual a 0.7 metros, o agente chuta levemente a bola e anda de forma a conduzi-la ao gol direito.
- **Chutar bola para gol direito:** Caso a distância até a bola seja menor ou igual a 0.7 metros, o agente chuta a bola na direção do gol direito.
- **Carregar bola para o gol esquerdo:** Caso a distância até a bola seja menor ou igual a 0.7 metros, o agente chuta levemente a bola e anda de forma a conduzi-la ao gol esquerdo.
- **Chutar bola para gol esquerdo:** Caso a distância até a bola seja menor ou igual a 0.7 metros, o agente chuta a bola na direção do gol esquerdo.

Nos Códigos 3.5, 3.6 e 3.7 são descritos os algoritmos dos comportamentos de forma geral. Com poucas adaptações chega-se aos comportamentos específicos para direita ou esquerda do campo.

```
função AproximarBola(s):  
  se s.BolaVisível:  
    ângulo da bola <- s.DireçãoBola + s.DireçãoPescoço  
    comando <- Arrancar(60, ângulo da bola)  
  retorna comando
```

Código 3.5: Algoritmo do comportamento de se aproximar da bola

```
função CarregarBola(s):  
  se s.DistânciaBola < 0.7:  
    ângulo <- -s.Orientação  
    comando <- Chutar(5, ângulo)  
  retorna comando
```

Código 3.6: Algoritmo geral do comportamento de carregar bola

```

função CarregarBola(s):
  Xgol, Ygol <- Gol.Posição
  se s.DistânciaBola < 0.7:
    ângulo <- atan2(s.Y - Ygol, s.X - Xgol) - s.Orientação
    comando <- Chutar(60, ângulo)
  retorna comando

```

Código 3.7: Algoritmo geral do comportamento de chutar bola para o gol

3.5.3 Parâmetros

- **Fator de desconto (γ):** Apesar do ambiente ser episódico, foi utilizado um fator de desconto de 0.99 devido ao fato de que a condição de término do episódio (fim de jogo) não ser observável através da discretização do estado utilizada.
- **Fator de aprendizagem (α):** O fator de aprendizagem foi definido inicialmente como 0.1 e foi reduzido exponencialmente multiplicando-o por 0.99999 ao final de cada episódio.
- **Fator de exploração (ϵ):** Para incentivar a exploração das possibilidades, o fator de exploração foi definido inicialmente como 0.9 e reduzido exponencialmente multiplicando-o por 0.99990 ao final de cada episódio. A cada ação tomada, o agente tem probabilidade ϵ de escolher uma ação aleatória. Além disso, para favorecer a exploração, no início de cada episódio era também sorteada uma posição inicial para o agente em seu lado do campo.

3.5.4 Recompensa

A recompensa foi definida em 3 partes a fim de guiar o agente na direção do aprendizado desejado. Todas as medições de recompensa foram feitas utilizando dados da simulação e não da percepção do agente.

- **Proximidade da bola R_1 :** Para que o agente tenha tendência a se aproximar da bola, foi definida uma recompensa negativa proporcional à distância d do agente em relação à bola, ou seja: $R_1 = -d * 0.001/6000$.
- **Velocidade da bola R_2 :** Para que o agente adquira o comportamento de chutar a bola em direção ao gol adversário, foi definida uma recompensa positiva proporcional à velocidade instantânea da bola em X (v_x), ou seja: $R_2 = v_x/6000$.
- **Gol R_3 :** Por fim, para incentivar que o agente fizesse gols, foi definida uma recompensa esparsa de valor 1 para cada gol realizado e -1 para gols contra, ou seja:

$$R_3 = \begin{cases} 1 & \text{se foi feito um gol no ciclo;} \\ -1 & \text{se foi feito um gol contra no ciclo;} \\ 0 & \text{caso não houver gol no ciclo.} \end{cases}$$

Deste modo, a cada instante de tempo foi atribuída uma recompensa $R = R_1 + R_2 + R_3$.

3.6 *Q-Learning* Duplo Tabular e Comportamentos Pré-Programados

Além da solução através de um aproximador de funções, foi testada uma solução utilizando o método tabular com os mesmos comportamentos e uma discretização dos estados descrita abaixo.

A codificação dos comportamentos, parâmetros α e γ e a recompensa foram mantidos como descrito nas subseções 3.5.2, 3.5.3 e 3.5.4 respectivamente. O parâmetro ϵ foi mantido em 0.9 porém com um decaimento mais lento, sendo multiplicado por 0.99996 ao final de cada episódio.

3.6.1 Codificação dos Estados

O estado percebido pelo agente é dado pela combinação dos seguintes fatores:

- **Distância até a bola.** A distância D até a bola foi discretizada de acordo com a seguinte função:

$$\begin{cases} 0 & \text{se } D < 0.7; \\ \lfloor \log_2(\frac{D}{0.7}) \rfloor & \text{se } 0.7 \leq D \text{ e } D < 0.7 \times 2^6; \\ 6 & \text{se } D \geq 0.7 \times 2^6; \end{cases} \quad (3.2)$$

ou seja, a distância percebida até a bola varia entre 0 e 6 com resolução cada vez menor à medida que o agente se afasta da bola. O fator 0.7 foi inserido na função devido ao fato de que esta é a distância mínima que permite que o agente chute a bola.

Caso o jogador possa enxergar a bola, a distância D é recebida diretamente do sensor. Caso contrário, a distância D é estimada com base na última posição percebida da bola.

- **Direção da bola:** A direção da bola foi dividida em 24 setores de 15° cada. O ângulo de visão do jogador é de $\pm 30^\circ$. Caso a bola não esteja visível, a direção da bola é estimada com base na última posição percebida.
- **Posição do jogador em X:** A posição estimada do jogador em X foi discretizada em 10 intervalos de tamanho 11.5.
- **Posição do jogador em Y:** A posição estimada do jogador em Y foi discretizada em 7 intervalos de tamanho aproximado 11.14.
- **Direção do jogador:** A direção estimada do jogador em relação ao eixo horizontal foi discretizada em 24 setores de 15° cada.

Com isso, o número total de estados possíveis é dado pelo produtório da quantidade de possibilidades em cada um dos itens acima totalizando 282240 estados.

3.7 *Q-Learning* Duplo Tabular e Ações Elementares

Após o resultado positivo obtido utilizando comportamentos pré-programados, foi testada uma abordagem *end-to-end* utilizando a discretização de estados descrita na subseção 3.6.1 e a

recompensa descrita na subseção 3.5.4.

3.7.1 Codificação das Ações

Para simplificar o vasto espaço de ações disponíveis, foi selecionado um conjunto discreto de 13 ações:

- **Ação nula:** O agente apenas espera até o próximo ciclo.
- **Virar-se:** O agente tem a opção de virar-se 7° , 15° ou 31° para ambas as direções, totalizando 6 ações de rotação possíveis.
- **Correr:** É possível correr em frente (0°) ou a 30° em ambas as direções, sempre com potência 50, totalizando 3 ações de corrida possíveis.
- **Chutar:** Caso a distância até a bola seja menor ou igual a 0.7 metros, o jogador tem a opção de chutá-la para frente ou em um ângulo de 45° para esquerda ou direita, totalizando 3 ações de chute possíveis. Caso a bola não esteja próxima o suficiente, nada acontece.

3.7.2 Parâmetros

- **Fator de desconto (γ):** Foi mantido tal como descrito na subseção 3.5.3.
- **Fator de aprendizagem (α):** O fator de aprendizagem foi definido como $\frac{6000}{6000+N}$ sendo N o número de episódios experienciados pelo jogador.
- **Fator de exploração (ϵ):** O fator de aprendizagem foi definido como $\frac{8000}{8000+N}$ sendo N o número de episódios experienciados pelo jogador.

No capítulo seguinte, é feita a análise dos resultados dos experimentos descritos neste capítulo.

Capítulo 4

Resultados Experimentais

4.1 Sarsa com Aproximador de Funções e Comportamentos Pré-Programados

O treinamento com o algoritmo Sarsa contou com 30.000 episódios com a utilização da rede neural multicamadas descrita na Subseção 3.5 e salvando os retornos obtidos em cada episódio.

O gráfico da Figura 4.1 mostra em conjunto o retorno a cada episódio e a média móvel do retorno com janela de 1000 treinamentos de partidas. Como a recompensa definida apresenta parcelas R_1 e R_2 com ordem de grandeza muito menores que R_3 , que corresponde ao gol contra ou a favor, o retorno pode ser interpretado como aproximadamente o número de gols realizados durante o episódio.

É possível observar que há uma tendência de subida do retorno com o passar dos episódios experienciados. Porém, em grande parte dos episódios, o agente não realizou nenhum gol, refletido pelo baixo valor da média de 1000 episódios, que precisou ser multiplicado por 10 para melhor visualização gráfica.

É interessante ressaltar o alto custo computacional deste tipo de treinamento devido à utilização de redes neurais. Dessa forma, a quantidade de amostras possíveis de serem coletadas em tempo hábil foi drasticamente reduzida.

Além disso, a utilização de métodos aproximados apresenta um problema de duplo aprendizado: deseja-se aprender a política ótima enquanto se aprende a aproximar esta política ótima desconhecida por meio de uma rede neural. Esse fator contribui negativamente no tempo para convergência da política. Em métodos tabulares, apesar do maior custo de memória, esse problema é inexistente.

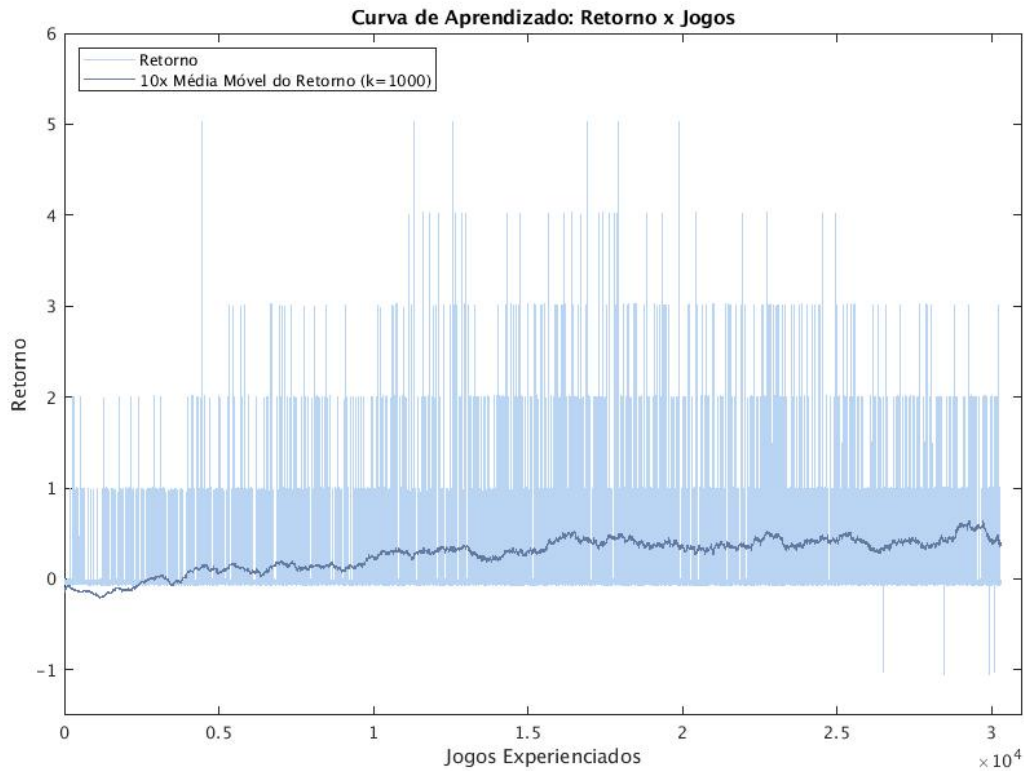


Figura 4.1: Curva de aprendizado do agente com comportamentos pré-programados utilizando Sarsa aproximado.

4.2 *Q-Learning* Duplo Tabular e Comportamentos Pré-Programados

Para este experimento, o aproximador de função foi substituído por uma tabela e o algoritmo Sarsa pelo Q-Learning. Foram executados 3 treinamentos de 100.000 episódios e a tabela Q completa e o histórico dos retornos obtidos pelo agente ao longo do treinamento foram salvos.

O gráfico da Figura 4.2 mostra o histórico médio dos 3 treinamentos. Observa-se que há uma estagnação do retorno por volta das 60.000 amostras, o que pode indicar a necessidade de ajuste no decaimento do fator de exploração para que o agente explore novas possibilidades por mais episódios. Isso mostra a importância de buscar equilíbrio no dilema da exploração e aproveitamento, discutido na Seção 2.2.5.

A estagnação pode indicar, também, existência de um limite superior para o desempenho do agente devido à menor flexibilidade da política aprendida, ou seja, o agente só é capaz de construir a política a partir dos comportamentos pré-programados.

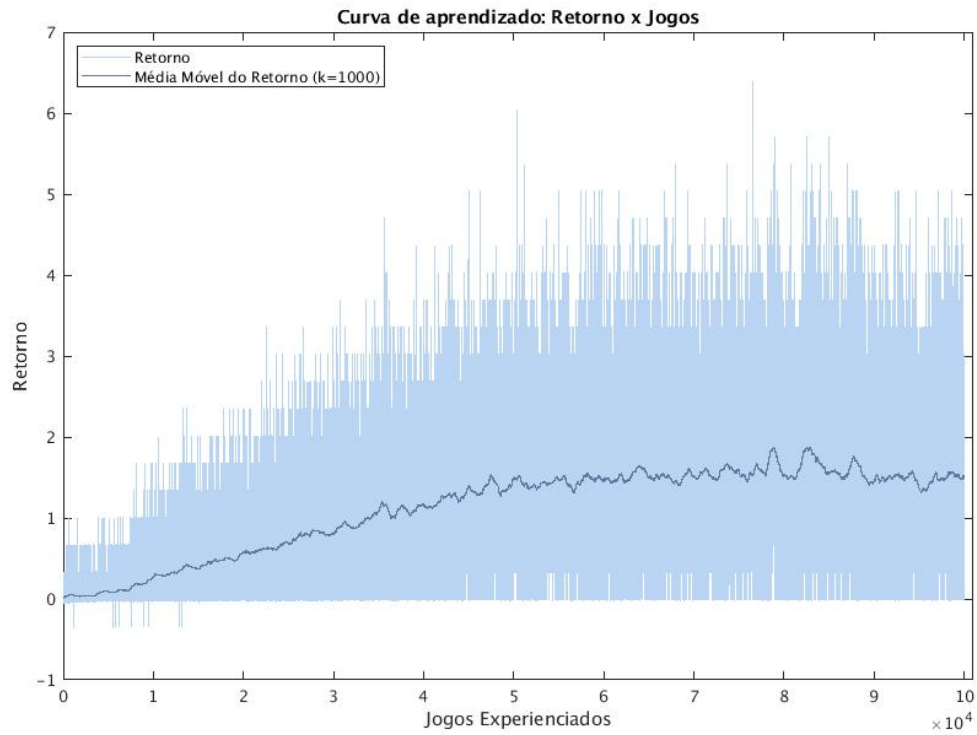


Figura 4.2: Curva de aprendizado do agente com comportamentos pré-programados.

A Figura 4.3 ilustra o agente único treinado conduzindo a bola e realizando um gol conforme a política aprendida.

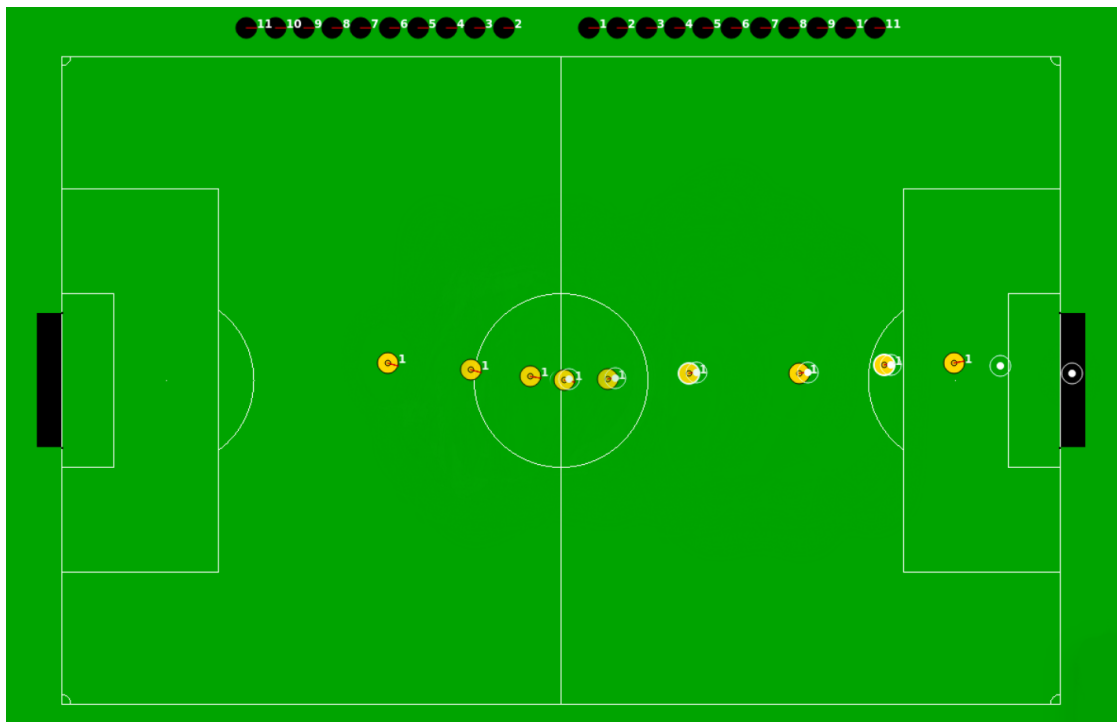


Figura 4.3: Sobreposição de sequência de imagens do agente único fazendo gol.

A Figura 4.4 mostra o histórico de retornos para um treinamento mais longo, de 200.000 episódios. Nela percebe-se que, após a cessação da exploração, o agente estabiliza seu desempenho um pouco abaixo do máximo obtido.

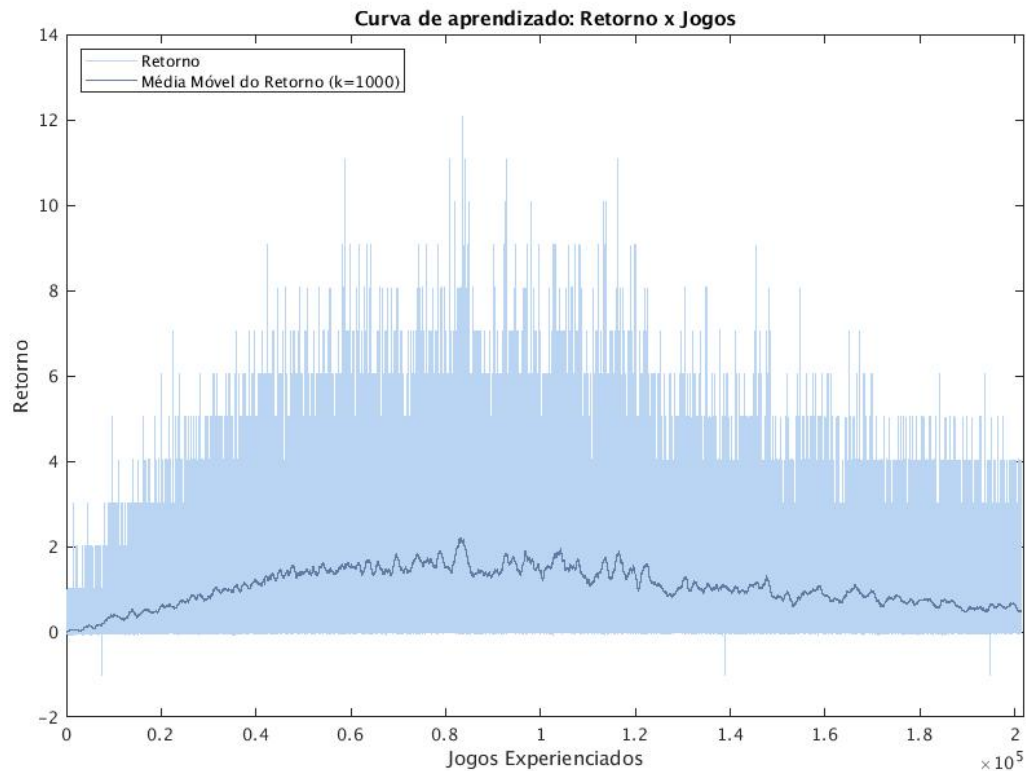


Figura 4.4: Curva de aprendizado do agente com comportamentos pré-programados para treinamento longo.

Utilizando os pesos obtidos em uma das execuções do treinamento é possível observar cenários variados da política aprendida com um único agente e a bola em campo. A Figura 4.5 ilustra quatro situações distintas construídas com sobreposição de imagens a cada 50 ciclos. Em três delas, o agente marcou um gol em tempos variados, ora tocando a bola até a área do gol, ora chutando de uma distância maior. A quarta situação, no canto inferior direito, ilustra uma situação em que o agente não realizou gols e ficou parado após tocar a bola para trás.

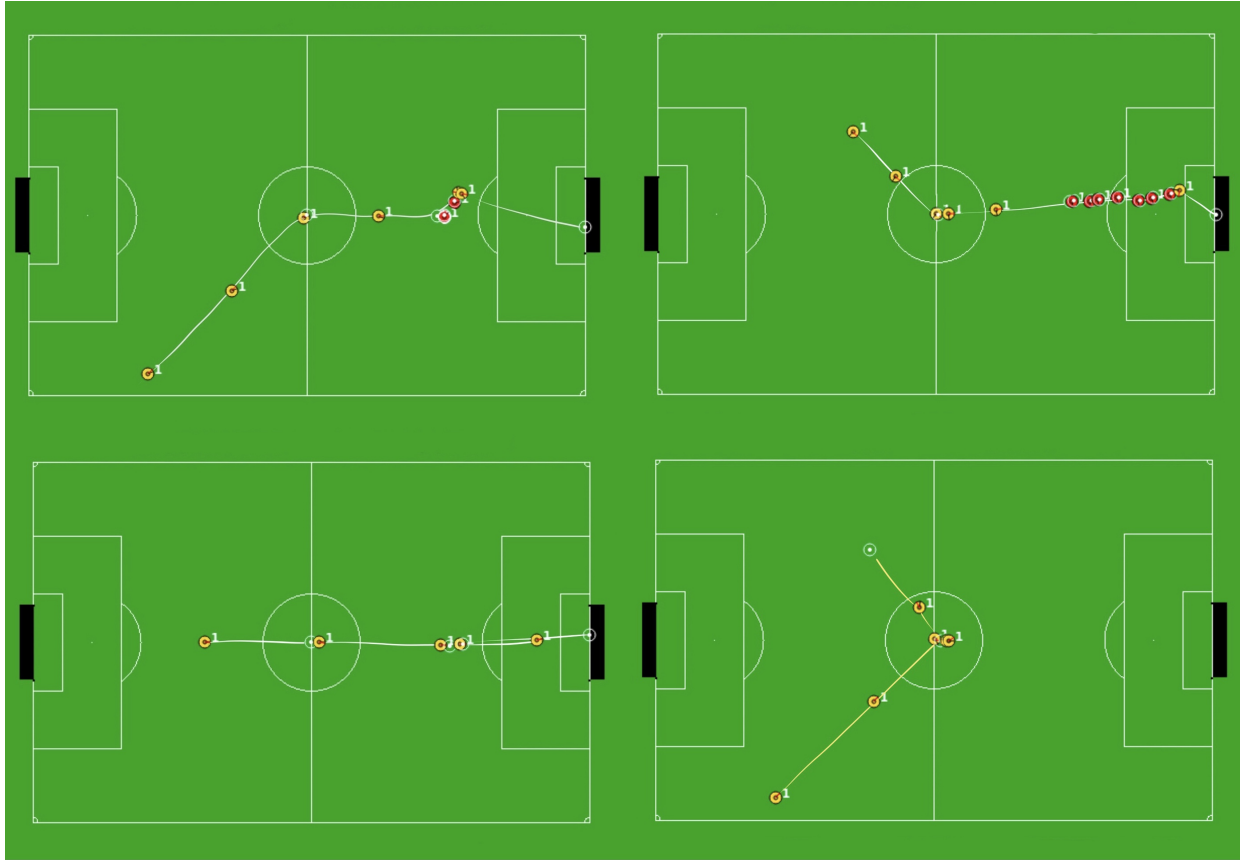


Figura 4.5: Montagem ilustrando 4 exemplos da política aprendida através de sobreposição de imagens.

4.3 *Q-Learning* Duplo Tabular e Ações Elementares

Substituindo os comportamentos pré-programados por uma seleção de ações elementares, foram executados 3 treinamentos distintos de 100.000 treinamentos de partidas a fim de suavizar o elemento sorte nos resultados. Após cada um dos treinamentos, foram salvos a tabela Q completa e o histórico dos retornos obtidos pelo agente ao longo do treinamento.

A Figura 4.6 mostra esse histórico. É interessante observar que, com o decaimento dos fatores de exploração e de aprendizagem, após 100.000 episódios, ambos eram $\epsilon \approx 0.074074$ e $\alpha \approx 0.056604$, ou seja, o agente já executava na maior parte dos ciclos a política aprendida. Para cada jogo foi feita a média entre os 3 retornos observados em cada um dos treinamentos.

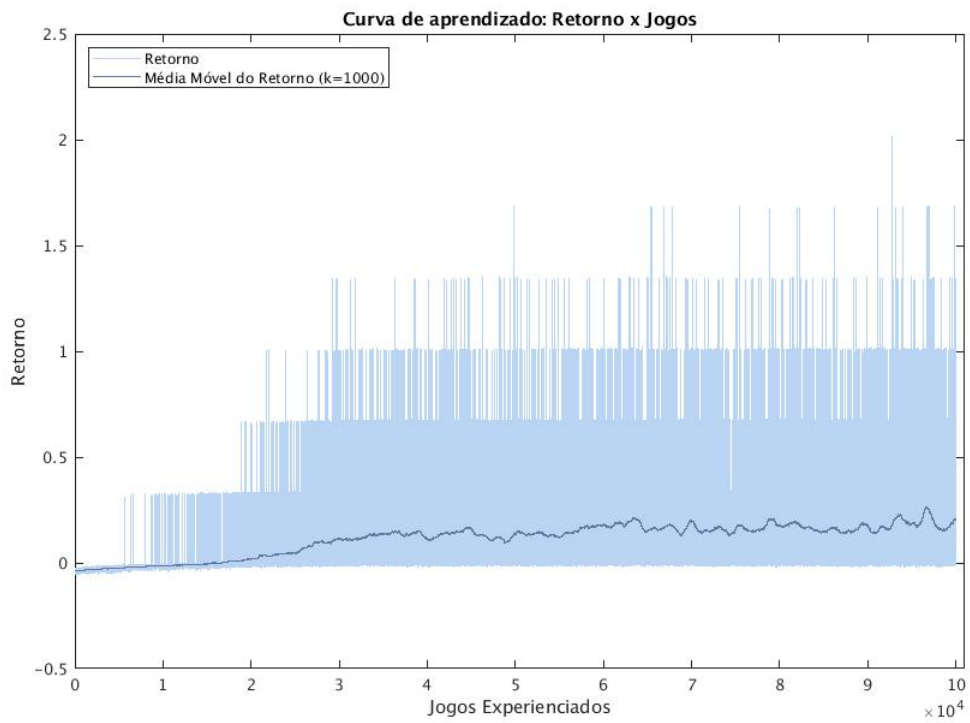


Figura 4.6: Curva de aprendizado do agente com ações elementares.

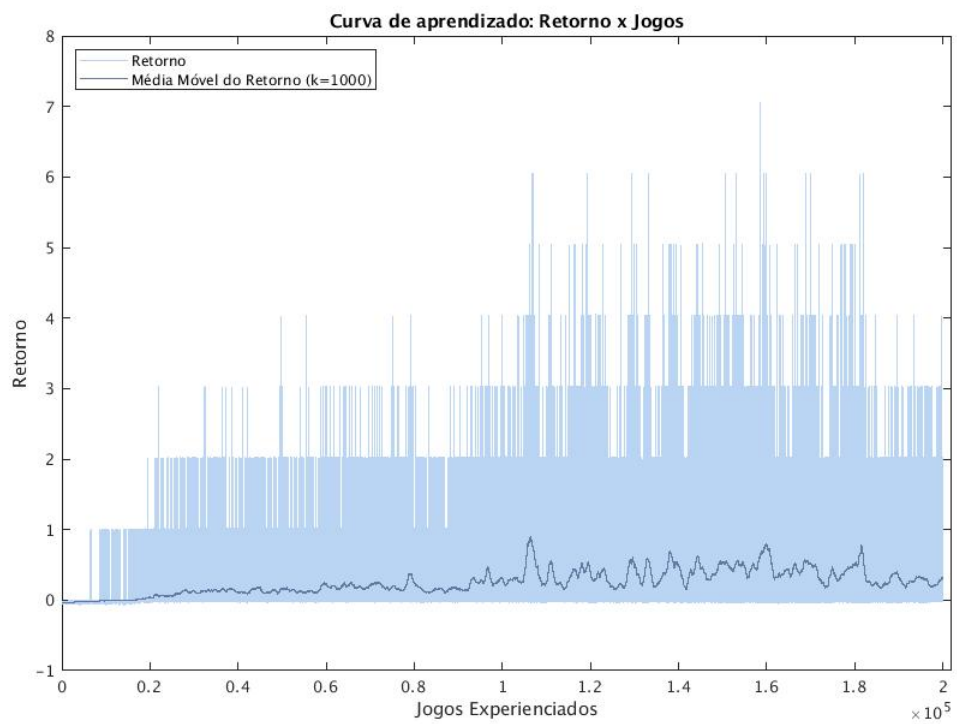


Figura 4.7: Curva de aprendizado para treinamento longo.

Além disso, assim como na Seção 4.2, foi executado um treinamento de 200.000 episódios a fim de observar a aprendizagem por um período mais longo. Na Figura 4.7, observa-se que o agente continua melhorando seu desempenho até próximo do fim do treinamento, o que pode ser indício de que o aprendizado com ações elementares de fato permite mais flexibilidade na política aprendida.

Apesar da grande quantidade de experiência a que o agente teve acesso, nota-se na Figura 4.7 que o crescimento de seu desempenho é bastante limitado, sequer atingindo a média de 1 gol por episódio. Isso é um indicativo do altíssimo custo computacional de soluções utilizando ações elementares como a utilizada no experimento, especialmente quando o espaço de estados é vasto.

O capítulo seguinte relata as conclusões acerca desse trabalho, sua contribuição para a área e discorre sobre possíveis trabalhos para continuação do desenvolvimento da plataforma, especialmente dentro do contexto da Universidade de Brasília.

Capítulo 5

Conclusões

Este projeto teve como objetivo a implementação de uma plataforma para desenvolvimento de aprendizagem por reforço em futebol de robôs, mais especificamente na categoria RoboCup Soccer Simulation 2D, e a utilização de técnicas de RL para realizar treinamentos de seleção de comportamentos e ações elementares para maximizar o número de gols feitos por um agente.

Ao pesquisar sobre a comunidade e equipes participantes das edições nacionais da competição, notou-se que a biblioteca de interfaceamento com o servidor *librsc* e o time base *agent2d* - ambos desenvolvidos no Japão por acadêmicos da equipe *HELIOS* - são amplamente utilizados. Apesar de ser uma plataforma bastante completa, sua dominância acaba limitando a inovação na categoria, mantendo as equipes dentro da arquitetura de solução proposta no *agent2d*.

A plataforma desenvolvida no decorrer deste trabalho ajuda a modernizar e diversificar a base de código utilizada pelas equipes. A plataforma desenvolvida encontra-se disponível em repositórios do GitHub: *rcssggb/ggb-lib*, para a biblioteca de comunicação, e *rcssggb/ggb-single*, onde se encontra a rotina de treinamento. Ambas podem ser utilizadas livremente. A plataforma se conecta via protocolo UDP ao servidor que executa a partida e abstrai a codificação e decodificação de mensagens enviadas ao servidor. Algoritmos de RL implementados no decorrer do trabalho validaram o funcionamento da plataforma.

Foram realizados três treinamentos distintos e seus resultados foram analisados conforme o retorno alcançado pelo agente. Inicialmente, o algoritmo Sarsa com comportamentos pré-programados e aproximador de funções mostrou um retorno crescente com o passar dos episódios, porém o custo computacional do treinamento mostrou-se muito alto, tornando difícil a testagem de soluções desse tipo.

A fim de aumentar o número de amostras, o estado foi discretizado e implementou-se o algoritmo *Q-Learning* duplo com os mesmos comportamentos e utilizando um método tabular. O número de amostras possíveis de serem coletadas aumentou consideravelmente e o treinamento se mostrou mais estável, além de obter retornos mais altos.

Por fim, foi realizado um treinamento com ações elementares, permitindo que o agente aprendesse estratégias por si só. O treinamento demonstrou retornos menores que o anterior, porém a

aprendizagem pareceu continuar por mais tempo, o que sugere que abordagens desse tipo sejam capazes de aproveitar melhor quantidades maiores de amostras.

5.1 Trabalhos Futuros

A plataforma desenvolvida abre muitas possibilidades para pesquisa e desenvolvimento dentro do contexto da RoboCup Soccer Simulation 2D. É possível estudar diversas técnicas de RL a fim de cumprir a proposta da competição: desenvolver um time de futebol completo com agentes capazes de cooperar entre si e disputar contra outras equipes participantes.

Apesar da abordagem simplificada feita neste trabalho, o formato oficial da competição envolve um ambiente multiagente e adversarial que introduz diversos desafios ao aprendizado uma vez que o ambiente se torna mais dinâmico e, com isso, mais imprevisível.

Além disso, a família de técnicas de RL ideais para um cenário como a RCSS é a família métodos de gradientes de políticas, ou seja, o agente busca estimar a política em si e não a função de valor de ação Q . Com essas técnicas, é possível parametrizar a política de modo a cobrir todo o espaço de ações e não apenas um subconjunto discreto dele.

Por fim, em um cenário competitivo é importante manter a plataforma atualizada e capaz de abstrair todas as funcionalidades do servidor. Uma das funcionalidades faltantes, por exemplo, é a da investida, também conhecida coloquialmente como "carrinho".

As equipes acadêmicas são um ambiente propício para o aprendizado prático de diversas técnicas aprendidas durante a graduação e a Universidade de Brasília tem uma longa tradição de equipes participantes em diversas categorias de robótica, incluindo outras categorias de futebol de robôs. A organização de uma equipe para desenvolvimento utilizando a plataforma a fim de participar da RCSS criaria diversas oportunidades de aprendizado e pesquisa.

REFERÊNCIAS BIBLIOGRÁFICAS

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press, 2nd ed., 2018.
- [2] M. Chen, K. Dorer, E. Foroughi, F. Heintz, Z. Huang, S. Kapetanakis, K. Kostiadis, J. Kumeneye, J. Murray, I. Noda, O. Obst, P. Riley, T. Steffens, Y. Wang, and X. Yin, “Users Manual - RoboCup Soccer Server (for Soccer Server version 7.07 and later),” 2003.
- [3] M. Guarnieri, “The roots of automation before mechatronics [historical],” *IEEE Industrial Electronics Magazine*, vol. 4, no. 2, pp. 42–43, 2010.
- [4] N. Wiener, *Cybernetics or Control and Communication in the Animal and the Machine*. Technology Press, 1948.
- [5] P. McCorduck, *Machines Who Think: A Personal Inquiry into the History and Prospects of Artificial Intelligence*. AK Peters Ltd, 2004.
- [6] B. Siciliano and O. Khatib, *Springer Handbook of Robotics*. Berlin, Heidelberg: Springer-Verlag, 2007.
- [7] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, 2016.
- [8] OpenAI, “OpenAI Five.” <https://blog.openai.com/openai-five/>, 2018.
- [9] A. K. Mackworth, “On seeing robots,” in *Computer Vision: Systems, Theory and Applications*, pp. 1–13, World Scientific, 1993.
- [10] H. Kitano, M. Asada, Y. Kuniyoshi, I. Noda, and E. Osawa, “Robocup: The robot world cup initiative,” in *Proceedings of the first international conference on Autonomous agents*, pp. 340–347, 1997.
- [11] A. Bai, F. Wu, and X. Chen, “Online planning for large markov decision processes with hierarchical decomposition,” *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 6, no. 4, p. 45, 2015.

- [12] T. Nakashima, H. Akiyama, Y. Suzuki, A. Ohori, and T. Fukushima, “Helios2018: Team description paper,” in *RoboCup 2018 Symposium and Competitions: Team Description Papers, Montreal, Canada*, 2018.
- [13] F. Mello, L. Ramos, M. Maximo, R. Ferreira, and V. Moura, “ITAndroids 2D soccer simulation team description 2012,”
- [14] M. R. Maximo, F. V. Coimbra, and H. F. Feitosa, “ITAndroids 2D soccer simulation team description 2019,”