

Dissertação apresentada à Pró-Reitoria de Pós-Graduação e Pesquisa do Instituto Tecnológico de Aeronáutica, como parte dos requisitos para obtenção do título de Mestre em Ciências no Programa de Pós-Graduação em Engenharia Aeronáutica e Mecânica, Área de Sistemas Aeroespaciais e Mecatrônica.

**Maria das Graças Silva**

**UMA ABORDAGEM SOBRE O DILEMA DO  
CUPIM FRENTE AO CONCRETO ARMADO  
UTILIZANDO DIFERENTES COMPOSIÇÕES  
CIMENTÍSTICAS**

Dissertação aprovada em sua versão final pelos abaixo assinados:

Prof. Dr. Adalberto Santos Dupont

Orientador

Prof<sup>a</sup>. Dr<sup>a</sup>. Doralice Serra

Coorientadora

Prof. Dr. John von Neumann

Pró-Reitor de Pós-Graduação e Pesquisa

Campo Montenegro  
São José dos Campos, SP - Brasil  
2015

**Dados Internacionais de Catalogação-na-Publicação (CIP)**  
**Divisão de Informação e Documentação**

Silva, Maria das Graças

Uma Abordagem Sobre o Dilema do Cupim Frente ao Concreto Armado Utilizando Diferentes Composições Cimentísticas / Maria das Graças Silva.

São José dos Campos, 2015.

34f.

Dissertação de Mestrado – Curso de Engenharia Aeronáutica e Mecânica. Área de Sistemas Aeroespaciais e Mecatrônica – Instituto Tecnológico de Aeronáutica, 2015. Orientador: Prof. Dr. Adalberto Santos Dupont. Coorientadora: Prof<sup>a</sup>. Dr<sup>a</sup>. Doralice Serra.

1. Cupim. 2. Dilema. 3. Construção. I. Instituto Tecnológico de Aeronáutica. II. Título.

**REFERÊNCIA BIBLIOGRÁFICA**

SILVA, Maria das Graças. **Uma Abordagem Sobre o Dilema do Cupim Frente ao Concreto Armado Utilizando Diferentes Composições Cimentísticas**. 2015. 34f. Dissertação de Mestrado – Instituto Tecnológico de Aeronáutica, São José dos Campos.

**CESSÃO DE DIREITOS**

NOME DA AUTORA: Maria das Graças Silva

TÍTULO DO TRABALHO: Uma Abordagem Sobre o Dilema do Cupim Frente ao Concreto Armado Utilizando Diferentes Composições Cimentísticas.

TIPO DO TRABALHO/ANO: Dissertação / 2015

É concedida ao Instituto Tecnológico de Aeronáutica permissão para reproduzir cópias desta dissertação e para emprestar ou vender cópias somente para propósitos acadêmicos e científicos. A autora reserva outros direitos de publicação e nenhuma parte desta dissertação pode ser reproduzida sem a autorização da autora.

---

Maria das Graças Silva  
Av. Cidade Jardim, 679  
12.233-066 – São José dos Campos–SP

# UMA ABORDAGEM SOBRE O DILEMA DO CUPIM FRENTE AO CONCRETO ARMADO UTILIZANDO DIFERENTES COMPOSIÇÕES CIMENTÍSTICAS

Maria das Graças Silva

Composição da Banca Examinadora:

Prof. Dr.	Alan Turing	Presidente	-	ITA
Prof. Dr.	Adalberto Santos Dupont	Orientador	-	ITA
Prof <sup>a</sup> . Dr <sup>a</sup> .	Doralice Serra	Coorientadora	-	OVNI
Prof. Dr.	Linus Torwald		-	UXXX
Prof. Dr.	Richard Stallman		-	UYYY
Prof. Dr.	Donald Duck		-	DYSNEY
Prof. Dr.	Mickey Mouse		-	DISNEY

ITA

Aos amigos da Graduação e Pós-Graduação do ITA por motivarem tanto a criação deste template pelo Fábio Fagundes Silveira quanto por motivarem a mim e outras pessoas a atualizarem e aprimorarem este excelente trabalho.

# Agradecimentos

Primeiramente, gostaria de agradecer ao Dr. Donald E. Knuth, por ter desenvolvido o T<sub>E</sub>X.

Ao Dr. Leslie Lamport, por ter criado o L<sup>A</sup>T<sub>E</sub>X, facilitando muito a utilização do T<sub>E</sub>X, e assim, eu não ter que usar o Word.

Ao Prof. Dr. Meu Orientador, pela orientação e confiança depositada na realização deste trabalho.

Ao Dr. Nelson D'Ávila, por emprestar seu nome a essa importante via de trânsito na cidade de São José dos Campos.

Ah, já estava esquecendo... agradeço também, mais uma vez ao T<sub>E</sub>X, por ele não possuir vírus de macro :-)

*"If I have seen farther than others,  
it is because I stood on the shoulders of giants."*

— SIR ISAAC NEWTON

# Resumo

Aqui começa o resumo do referido trabalho. Não tenho a menor idéia do que colocar aqui. Sendo assim, vou inventar. Lá vai: Este trabalho apresenta uma metodologia de controle de posição das juntas passivas de um manipulador subatuado de uma maneira subótima. O termo subatuado se refere ao fato de que nem todas as juntas ou graus de liberdade do sistema são equipados com atuadores, o que ocorre na prática devido a falhas ou como resultado de projeto. As juntas passivas de manipuladores desse tipo são indiretamente controladas pelo movimento das juntas ativas usando as características de acoplamento da dinâmica de manipuladores. A utilização de redundância de atuação das juntas ativas permite a minimização de alguns critérios, como consumo de energia, por exemplo. Apesar da estrutura cinemática de manipuladores subatuados ser idêntica a do totalmente atuado, em geral suas características dinâmicas diferem devido a presença de juntas passivas. Assim, apresentamos a modelagem dinâmica de um manipulador subatuado e o conceito de índice de acoplamento. Este índice é utilizado na sequência de controle ótimo do manipulador. A hipótese de que o número de juntas ativas seja maior que o número de passivas ( $n_a > n_p$ ) permite o controle ótimo das juntas passivas, uma vez que na etapa de controle destas há mais entradas (torques nos atuadores das juntas ativas), que elementos a controlar (posição das juntas passivas).

# Abstract

Soccer 2D is a simulated league, two teams of 11 intelligent autonomous agents play a soccer game in two dimensions. Each agent represents a player, receives limited information about the game situation and under them has to decide which action upon many to take. This simulation has it's own version of penalty and the aim of this work is to develop an optimized behavior for the goalkeeper under this situation utilizing Deep Reinforcement Learning techniques.



# Lista de Figuras

FIGURA 1.1 – Ke Jie faces AlphaGo . . . . .	16
FIGURA 1.2 – Ke Jie is ultimately defeated . . . . .	17
FIGURA 1.3 – Distribution of answers of 1,100 candidates from Deloitte’s survey about AI in the Enterprise on how the technology levered their bu- siness production in comparison to their competitors . . . . .	17
FIGURA 1.4 – How AI brings benefit according to the surveyed candidates from Deloitte report . . . . .	21
FIGURA 1.5 – Simulated Soccer 2D League . . . . .	22
FIGURA 2.1 – Cupim cibernético. . . . .	24
FIGURA A.1 –Uma figura que está no apêndice . . . . .	33

# Lista de Tabelas

TABELA 2.1 – Exemplo de uma Tabela . . . . .	23
--	----

# Lista de Abreviaturas e Siglas

CTq	computed torque
DC	direct current
EAR	Equação Algébrica de Riccati
GDL	graus de liberdade
ISR	interrupção de serviço e rotina
LMI	linear matrices inequalities
MIMO	multiple input multiple output
PD	proporcional derivativo
PID	proporcional integrativo derivativo
PTP	point to point
UARMII	Underactuated Robot Manipulator II
VSC	variable structure control

# Lista de Símbolos

$a$	Distância
$\mathbf{a}$	Vetor de distâncias
$\mathbf{e}_j$	Vetor unitário de dimensão $n$ e com o $j$ -ésimo componente igual a 1
$\mathbf{K}$	Matriz de rigidez
$m_1$	Massa do cumpim
$\delta_{k-k_f}$	Delta de Kronecker no instante $k_f$

# Sumário

1	INTRODUCTION . . . . .	15
1.1	Motivation . . . . .	15
1.1.1	AI Super Powers . . . . .	15
1.1.2	Deep Learning Revolution . . . . .	18
1.2	Contextualization . . . . .	19
1.2.1	Simulated Soccer 2D . . . . .	19
1.3	Objective . . . . .	19
1.4	Scope . . . . .	20
1.5	Organization of this work . . . . .	20
2	MODELAGEM DINÂMICA DE CUPINS CIBERNÉTICOS . . . . .	23
2.1	Modelagem no espaço das juntas . . . . .	23
3	DEEP REINFORCEMENT LEARNING . . . . .	25
3.1	Markov Decision Processes . . . . .	25
3.1.1	Important Concepts . . . . .	26
3.1.2	The Bellman Equations . . . . .	27
3.1.3	About Solving It . . . . .	28
3.1.4	The Problem With The Tabular Solution . . . . .	28
3.2	Neural Networks . . . . .	29
3.2.1	Neurons . . . . .	29
4	CONCLUSÃO . . . . .	30
	REFERÊNCIAS . . . . .	31

---

APÊNDICE A – TÓPICOS DE DILEMA LINEAR . . . . .	33
A.1 Uma Primeira Seção para o Apêndice . . . . .	33
ANEXO A – EXEMPLO DE UM PRIMEIRO ANEXO . . . . .	34
A.1 Uma Seção do Primeiro Anexo . . . . .	34

# 1 Introduction

## 1.1 Motivation

In an afternoon of 2017, humanity had an unlikely champion. Ke Jie, the number one player in GO, was to battle against one of the world's most intelligent machines, AlphaGo; a powerhouse of artificial intelligence backed by one of the world's top technology company: Google. Since the number of possible positions on a Go board exceeds the quantity of atoms in the universe, defeating the world champion was the biggest deed possible for AI researchers at the time. Romantics said it was impossible for a machine to do so as it lacked spiritual power only a human could possess, engineers would agree but by a different reason, they pointed out for the fact that there was just too many possibilities for a computer to evaluate; these "theories" fell apart though. On this day, what was witnessed was not a clash between two equal titans, AlphaGo completely dominated Ke Jie. Over the course of three marathon matches of more than three hours each, the top ranked player tried every single approach: conservative, aggressive, defensive and unpredictable. Nothing worked, there were no openings, not a single sign of hope.

### 1.1.1 AI Super Powers

As described by (LEE, 2018), this event is the equivalent of Sputnik for China. Less than two months after Ke Jie resigned his last game to AlphaGo, the central government of China issued an ambitious plan to build artificial intelligence capabilities. China designed a clear strategy, set benchmarks for progress by 2020 and 2025 aiming to be the center of global innovation in artificial intelligence. Private money answered that call, Chinese venture-capital investors had already responded to that call, unleashing record amounts into AI startups and making up 48 percent of all AI venture funding globally, surpassing the United States for the first time.

Less than two months after Ke Jie resigned his last game to AlphaGo, the Chinese central government issued an ambitious plan to build artificial intelligence capabilities. It called for greater funding, policy support, and national coordination for AI development.



FIGURA 1.1 – Ke Jie faces AlphaGo

It set clear benchmarks for progress by 2020 and 2025, projected that by 2030 China would become the center of global innovation in artificial intelligence, leading in theory, technology, and application. By 2017, Chinese venture-capital investors had already responded to that call, pouring record sums into artificial intelligence startups and making up 48 percent of all AI venture funding globally, surpassing the United States for the first time.

The fact that the current global economic power and the one aspiring to be are massively investing on AI technology points at first hand that the topic is worth studying for. The State of AI in the Enterprise 2nd Edition conducted by Deloitte (LOUCKS DAVID SCHATSKY, 2018) gives some numbers on how AI is making the industry grow: “To obtain a cross-industry view of how organizations are adopting and benefiting from cognitive computing/AI, Deloitte surveyed 1,100 IT and line-of-business executives from US-based companies in Q3 2018. All respondents were required to be knowledgeable about their company’s use of cognitive technologies/artificial intelligence, and 90 percent have direct involvement with their company’s AI strategy.”. According to it, 82 per cent of the survey participants claim a positive financial return on their AI investment and 88 per cent plan to increase spending in 2019; 79 per cent agree that AI technologies empower people to make better decisions and 73 per cent believe AI will increase job satisfaction. Furthermore, the graphic 1.3 shows that around 75 per cent reported positive results when comparison to direct competitors was made.

What are the results AI can bring to justify this current fever? (AGRAWAL; GOLDFARB, 2018) puts it simply, AI is a technology that optimizes and facilitates decisions and





FIGURA 1.2 – Ke Jie is ultimately defeated

**AI helps organizations keep up with the (Dow) Joneses**

Relative to competitors, respondents say their company's adoption of AI has allowed them to . . .



Source: Deloitte State of AI in the Enterprise, 2nd Edition, 2018.

FIGURA 1.3 – Distribution of answers of 1,100 candidates from Deloitte’s survey about AI in the Enterprise on how the technology levered their business production in comparison to their competitors

predictions. For example, the article (LENTINO, 2019) talks about the Shanghai-based YITU Technology that has gained wide recognition for its Dragonfly Eye System, a facial scanning platform that can identify a person from a database of at least 2 billion people in a matter of seconds. (ZYMERGEN, 2019) talks about Zymergen, a company that uses machine learning to navigate the genomic search space in order to guide scientists to the precise set of genetic changes required to engineer cells to make a product of interest or to do it more efficiently. The applications are too numerous to cite them all but these two examples represent how broad the technology can go and how the future is being passed by it. The graphic 1.4 is from (LOUCKS DAVID SCHATSKY, 2018) and it shows more specifically how AI generate

### 1.1.2 Deep Learning Revolution

Machine learning, despite being this new driver of the future, is lucky to have survived a tumultuous half-century of research. For each great promising new lead there were big voids too, when weak practical results led to investments retraction.

There were 2 camps in the field: the “rule-based” approach and the “neural networks” approach. Researchers of the first attempted to teach computers to think by encoding a series of logical rules (If Then and Else). This method worked well for simple and well defined problems for example a game of Sudoku but when the universe of options was too large it did not respond well. Their insight was too question experts on the problem at hand and try to embody their wisdom into the machine (Hence why this method is also called the “expert systems”).

The “neural network” approach goes a different route by trying to reconstruct the brain itself, taking inspiration on the one and only source of intelligence we could perceive. This approach mimics the neurological architecture of the biological brain and instead of trying to transfer knowledge by force like the first one, this approach just feeds the neural networks with a plenitude of data and lets it find patterns on it.

(LEE, 2018) gives a simple illustration on how the techniques differ, if it is desired for the computer to recognize the picture of a cat, the “rule based” approach would work by trying, for starters, to embed the machine with the idea that a cat has 2 triangles above a circle. In the “neural network” approach, the machine would be fed with tones of pictures of cats and non-cats for it to build the features of what is a cat on its own.

To make things short, promising results were achieved during 1950s and 1960 by artificial neural networks but its core methodology found 2 problems that could not be solved at the time, first it needed a lot of data and second a lot of computational power. It quickly went out of fashion then and AI entered into one of its first “winters” during the 1970s.

It was not just the increased computing force that came with time which reanimated neural networks, a huge chunk of credit must be given to the optimization in the algorithm itself discovered by leading researcher Geoffrey Hinton and his Convolution Neural Networks (CITAR PAPER ORIGINAL). This new empowered neural network-now rebranded as “deep learning”- could outperform older models at a variety of tasks. However, the real demonstration occurred in 2012, when a neural network built by Hinton’s team overpowered the competition in an international computer vision contest.

## 1.2 Contextualization

RoboCup is an international scientific community aiming to advance the state of the art of intelligent robots. Its mission is that a team of androids will be able to beat the human team champion of the World Cup until the year 2050 (CITAR SITE AQUÍ). To achieve this goal, there is actually a plenitude of tasks to be solved, if you actually enumerate what a professional player must know, the difficulty stacks up.

For example, to kick a ball, a thing we humans almost take for granted, there is a complex mixture of muscle movements and positioning that a robot would have to mimic. Of course when you think about soccer you think about kicking but there is way more: to carry the ball a certain distance, to pass, to tackle, to receive a pass; these are all complex set of movements that requires good hardware and software coordination.

Soccer is not a pure athletic sport though, that means victory is not only decided by being more able, stronger or faster. It is a versus sport, the players must decide how to act in different conditions, when to run, where to position themselves and what movement to make according to the enemy same actions.

### 1.2.1 Simulated Soccer 2D

Directly from the RoboCup official site (CITAR SITE) : “In the 2D Simulation League, two teams of eleven autonomous software programs (called agents) each play soccer in a two-dimensional virtual soccer stadium represented by a central server, called SoccerServer. This server knows everything about the game, i.e. the current position of all players and the ball, the physics and so on. The game further relies on the communication between the server and each agent. On the one hand each player receives relative and noisy input of his virtual sensors (visual, acoustic and physical) and may on the other hand perform some basic commands (like dashing, turning or kicking) in order to influence its environment.”

In other words, this category focus on the behavior of a soccer player and disregard body-physics implementations.

## 1.3 Objective

The objective of this bachelor’s thesis is to use Deep Learning in order to optimize the behavior of the goalkeeper agent in Simulated Soccer 2D, under the more constricted penalty scenario in which the goalkeeper faces in a one-on-one dispute against an attacker that starts with the ball and is fairly distanced from the goalkeeper whose initial posi-

tion lies in the goal area. Currently, unlike real soccer game, the rules accept that the goalkeeper can move freely during the penalty.

## 1.4 Scope

In this work, Reinforcement Learning algorithms applied to models based on Deep Neural Networks will be approached in order to find, through gradient-based optimization techniques, optimal policies for the behavior of the goalkeeper. The initial intent is to use the algorithm known as PPO (Proximal Policy Optimization).

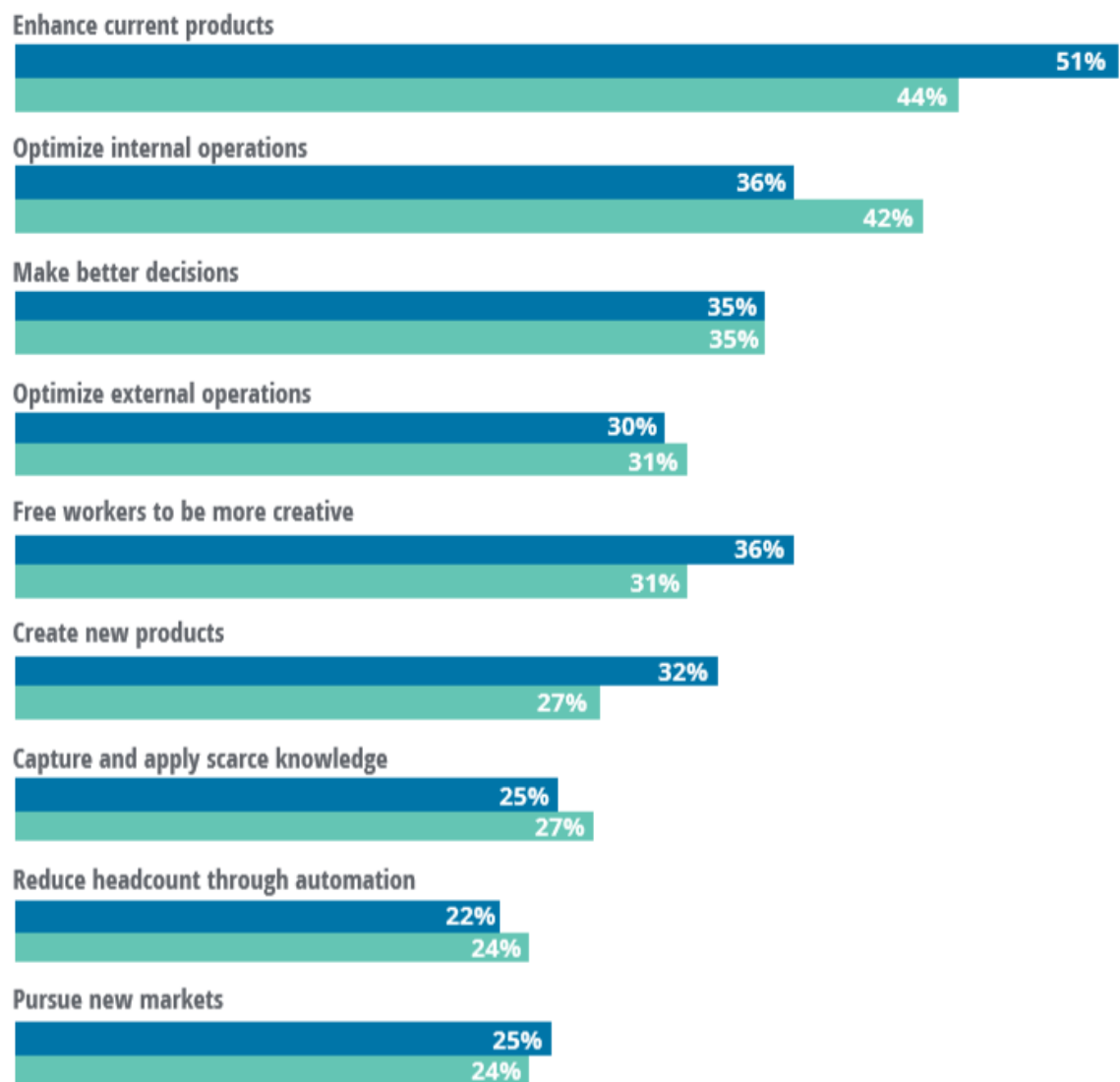
## 1.5 Organization of this work

Escrever depois

## AI's leading benefits are enhanced products and processes—and better decisions

Respondents rating each a top-three AI benefit for their company

■ 2017 ■ 2018



Source: Deloitte State of AI in the Enterprise, 2nd Edition, 2018.

FIGURA 1.4 – How AI brings benefit according to the surveyed candidates from Deloitte report



FIGURA 1.5 – Simulated Soccer 2D League

## 2 Modelagem Dinâmica de Cupins Cibernéticos

### 2.1 Modelagem no espaço das juntas

Manipuladores subatuados diferem dos totalmente atuados pois são equipados com um número de atuadores que é sempre menor que o número de graus de liberdade (GDL). Portanto, nem todos os GDL podem ser controlados ativamente ao mesmo tempo (SBORNIAN, 2004). Por exemplo, com um manipulador planar de 3 juntas equipado com dois atuadores, ou seja, duas juntas ativas e uma passiva, pode-se controlar ao mesmo tempo duas das juntas a qualquer instante, mas não todas. Para controlar todas as juntas de um manipulador subatuado, deve-se usar um controle sequencial. Este princípio foi provado pela primeira vez por arai usando argumentos dinâmicos linearizados (JOEA; JOHN, 2003), e é a base para a modelagem no espaço das juntas e no espaço Cartesiano. A Tabela 2.1 apresenta os resultados (ASSENMACHER *et al.*, 1993; ??; CAROMEL *et al.*, 1998).

Devido ao fato de que no máximo  $n_a$  coordenadas generalizadas (ângulos das juntas ou variáveis cartesianas) podem ser controladas num dado instante, o vetor de coordenadas generalizadas é dividido em duas partes, representando as coordenadas generalizadas ativas e as coordenadas generalizadas passivas (CALLAGHAN *et al.*, 1995).

Considerando um robô manipulador rígido, malha aberta, e de  $n$ -juntas em série. Seja  $q$  a representação de seu vetor de posição angular das juntas e  $\tau$  a representação de seu

TABELA 2.1 – Exemplo de uma Tabela

Parâmetro	Unidade	Valor da simulação	Valor experimental
Comprimento, $\alpha$	$m$	8, 23	8, 54
Altura, $\beta$	$m$	29, 1	28, 3
Velocidade, $v$	$m/s$	60, 2	67, 3



FIGURA 2.1 – Cupim cibernético.

vetor de torque. A equação dinâmica pelo método de Lagrange é dada por:

$$\frac{d}{dt}\left(\frac{\partial L}{\partial \dot{q}}\right) - \frac{\partial L}{\partial q} = \tau^T. \quad (2.1)$$

O Lagrangiano  $L$  é definido como a diferença entre as energias cinética e potencial do sistema:

$$L = T - P \quad (2.2)$$

A energia cinética total dos ligamentos é representada:

$$T = \frac{1}{2} \dot{q}^T M(q) \dot{q} \quad (2.3)$$



# 3 Deep Reinforcement Learning

## 3.1 Markov Decision Processes

A MDP (Markov Decision Process) is a mathematical model of the reinforcement learning problem. It is composed by an agent, the decision maker, and the thing it interacts with, the environment. These two interact continually, with the environment presenting a situation to the agent in the form of a state, and the agent selecting an action in regards to this state. This same environment also presents the agent with numerical rewards for each state accessed, the goal of the agent is to maximize the sum of these rewards by making good choices of actions.

In short, the MDP is defined by:

- A set of states  $s \in S$
- A set of actions  $a \in A$
- A transition function  $T(s, a, s')$ 
  - Probability that  $a$  from  $s$  leads to  $s'$ , i.e.,  $P(s' | s, a)$
  - Also called the model of the dynamics
- A reward function  $R(s)$
- A start state
- Maybe a terminal state

For each time step, the environment returns its current representation in the form of a state  $s$ . The agent then selects an action and the transition function  $T(s, a, s')$  returns the probability of reaching state  $s'$  from state  $s$  given the action chosen. On the next time step, the agents finds itself on a new state.

### 3.1.1 Important Concepts

#### 3.1.1.1 Reward Function

At each reached state, the agent accumulates a reward returned by a reward function  $R(s)$ . The purpose of this functions is to guide the agent, good states are scored with positive values, while bad states are scored with negative values. By making the agent act on the desire of maximizing his reward score, the agent will act optimally as it is intentioned.

#### 3.1.1.2 Transition Function

The MDP is not necessarily deterministic, after an agent chooses an action the next state is ruled by a probability distribution. The transaction function  $T(s, a, s')$  return the probability that an agent in state  $s$  performing action  $a$  will land in state  $s'$ .

#### 3.1.1.3 Discounting

In a MDP, it's reasonable to maximize the sum of rewards but it's also reasonable to make preference over rewards now than later. One solution is to make rewards decay exponentially in value. So for each episode, all rewards available are multiplied by a  $\gamma$  factor as illustrated by (FIGURA). This factor is called discount factor.

#### 3.1.1.4 Utility

Utility is, by definition, the sum of the discounted rewards for a given trajectory of the agent.

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (3.1)$$

In recursive form:

$$G_t \doteq R_{t+1} + \gamma G_{t+1} \quad (3.2)$$

#### 3.1.1.5 Policy

For MDPs, an *optimal policy*  $\pi^* : S \rightarrow A$  is desired. To put it simply, a *policy*  $\pi$  works as a map of instructions. Given a state  $s$ , the policy return which action  $a$  to take. Therefore, an optimal policy has an action for each state that maximizes the expected utility when followed. (FIGURA)

### 3.1.1.6 Value Function

The *value function* of a state  $s$  under a policy  $\pi$ , denoted  $v_\pi(s)$ , is the expected return when starting in  $s$  and following  $\pi$  thereafter. Formally:

$$v_\pi(s) \doteq \mathbf{E}_\pi[G_t \mid S_t = s] \quad (3.3)$$

where  $\mathbf{E}_\pi$  represents the expected value of a random variable given that the agent follows policy  $\pi$ .

### 3.1.1.7 Q-State

A Q-State is a tuple formed by a state  $s$  and an action  $a$ .

### 3.1.1.8 Action Value Function

Just like the Value Function, the *action value function* defines value for the Q-States. The value of taking action  $a$  in state  $s$  under the policy  $\pi$ , denoted  $q_\pi(s, a)$ , as the expected return starting from  $s$ , taking the action  $a$ , and thereafter following policy  $\pi$ . Formally:

$$q_\pi(s, a) \doteq \mathbf{E}_\pi[G_t \mid S_t = s, A_t = a] \quad (3.4)$$

## 3.1.2 The Bellman Equations

### 3.1.2.1 Optimal Quantities

- $V^*(s)$  is the expected utility starting in  $s$  and acting optimally.
- $Q^*(s, a)$  is the expected utility starting out having taken action  $a$  from state  $s$  and (thereafter) acting optimally.
- $\pi^*(s)$  returns the optimal action from state  $s$

### 3.1.2.2 The Equations

$$V^*(s) = \max_a Q^*(s, a) \quad (3.5)$$

$$Q^*(s, a) = \quad (3.6)$$

### 3.1.3 About Solving It

To make things clear, the solution required is the optimal policy, the set of instructions that maximizes the utility. If the transition function is known and well defined then the solution is easily found by means of dynamic programming, no reinforcement learning involved. Of course, in the real world, this transition function cannot be known a priori, that's where reinforcement learning comes.

The most basic reinforcement learning algorithms have the agent update state values in a series of trainings, multiple sessions of sampling substitute the transition function. The idea is to have the agent explore the environment by doing, at first, random actions but as time passes and almost all states are decently evaluated it progressively abandons the random nature of exploration in favor of a greedy behavior in which it chooses the action that will return the best benefit, ergo the best reward.

In Q-Learning algorithm, for example, the Q-States values are updated by the following expression 3.7:

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + (\alpha) \left[ r + \gamma \max_{a'} Q(s', a') \right] \quad (3.7)$$

where  $\alpha$  is a factor that averages the current value with the next value. The expression is very intuitive, the Q value is updated by an average between its current value and the reward the state gives plus the discounted maximum Q-Value accessible by the current state.

By doing a series of training, the Q-Values will converge and then the optimal policy may be extracted by using 3.8 below:

$$\pi^*(s) = \arg \max_a Q^*(s, a) \quad (3.8)$$

### 3.1.4 The Problem With The Tabular Solution

Q-Learning solves the MDP but its solution is brute force styled. That's because the algorithm contemplates all states and its possible actions without any scrutiny. For real world problems, that's not computationally feasible. For instance, Go has a number of possible situations greater than the number of atoms in the whole observable universe.

A great improvement is to make state approximations. For example, if you look at a PacMan game, each configuration on the board is a different state, move a ghost or PacMan just a bit and a whole new state is reached. The insight here is that some states are actually very close to each other. In example, all states that has PacMan moving close to a ghost (which results in PacMans death) are equally bad states to be avoided.

So there is actually no need to register all these similar states if you can extract crucial information from them and hereby decide if they are bad or not. What is wanted is that the machine, by looking at these relevant informations, learns what they mean in order to classify states.

The role of function approximation belongs to Neural Networks in Deep Learning.

## 3.2 Neural Networks

A Neural Network is a learning model whose goal is essentially function approximation. For example, for a classifier,  $y = f^*(x)$  maps an input  $x$  to a category  $y$ . A neural network defines a mapping  $y = f(x, \theta)$  and learns the value of the parameters  $\theta$  that result in best function approximation. They are called networks because they are commonly represented by composing together many different functions. The model is illustrated by a directed acyclic graph describing how the functions are composed together. For example, there could be three functions  $f_1$ ,  $f_2$  and  $f_3$  connected in a chain, to form  $f(x) = f_3(f_2(f_1(x)))$ .

### 3.2.1 Neurons

Neurons are the single unit from a Neural Network. They are an abstraction of a vector of *weights*  $\mathbf{w}$  and a number  $b$  called *bias*. Upon receiving an input  $\mathbf{x}$ , the neuron outputs:

$$z = \mathbf{w} \cdot \mathbf{x} + b \quad (3.9)$$

Actually, because the values of  $\mathbf{w}$  and  $b$  are supposed to be changed by a learning process, it is important that the output changes in linear form in relation to them so these changes may actuate smoothly. To achieve it, the *sigmoid function* is often used, the final output is therefore:

$$\text{Neuron output} = \sigma(z) = \frac{1}{1 + e^{-z}} \quad (3.10)$$

where  $z$  is extracted from 3.9.

### 3.2.2 The Network

In graph representation, a neuron is a vertex and each incoming edge is a weight in the weight vector  $\mathbf{w}$ . Networks are composed of layers of neurons, the first layer receives the initial input while the subsequent layers receive input from their predecessor layer. The last layer returns the final output.

Feedforward neural networks are called networks because they are typically represented by composing together many different functions. The model is associated with a directed acyclic graph describing how the functions are composed together. For example, we might have three functions  $f(1)$ ,  $f(2)$ , and  $f(3)$  connected in a chain, to form  $f(x) = f(3)(f(2)(f(1)(x)))$ . These chain structures are the most commonly used structures of neural networks. In this case,  $f(1)$  is called the first layer of the network,  $f(2)$  is called the second layer, and so on. The overall

## 4 Conclusão

Neste trabalho realizou-se o projeto de uma metodologia de controle subótimo redundante da junta passiva de um manipulador com três graus de liberdade instantaneamente. Para este propósito usou-se nas formulações o vetor gradiente de uma função escalar que estima o acoplamento entre a junta passiva e as ativas desse manipulador. Aqui a redundância foi usada da melhor maneira possível sem focalizar o efeito global. Portanto, este método deve ser denominado de *controle ótimo local por redundância*. A principal vantagem dessa formulação é a computação em tempo real, que é necessária para o controle do manipulador experimental. Além disso esse método pode ser usado com diferentes tipos de controladores, uma vez que as alterações são feitas nas equações dinâmicas do manipulador.

A consequência direta observada nessa formulação é a redução dos torques na fase de controle da junta passiva, e consequente redução da energia elétrica gasta. Isso ocorre devido ao fato de que ao longo da trajetória do manipulador o índice de acoplamento de torque tende a ser maximizado, e portanto, menor é o torque necessário nos atuadores para se conseguir o posicionamento da junta passiva do manipulador.

Outros resultados indiretos obtidos são: um movimento mais uniforme e suave do manipulador e um tempo de acomodação menor tanto no posicionamento da junta passiva quanto das ativas, conforme podemos observar nos gráficos de desempenho dos resultados apresentados. Isso ocorre porque a maximização do acoplamento entre as juntas facilita o controle. Assim ocorrem menos picos de torque, e como as juntas ativas tem “menos trabalho” para posicionar a passiva estas se movem menos na direção contrária ao movimento daquelas, diminuindo assim as velocidades alcançadas e os tempos de posicionamento.

Uma extensão deste trabalho pode ser a implementação de um *controle ótimo global por redundância* da junta passiva do manipulador. Para isto pode-se fazer o planejamento *off-line* da trajetória das juntas de modo a minimizar a energia consumida. Alguns estudos foram feitos nesse sentido, usando o Princípio Mínimo de Pontryagin, mas sem resultados satisfatórios até o momento.

# Referências

AGRAWAL, J. G. A.; GOLDFARB, A. **Prediction Machines**. 1st. ed. Boston: Harvard Business Review Press, 2018.

ASSENMACHER, H.; BREITBACH, T.; BUHLER, P.; HÜBSCH, V.; SCHWARZ, R. Panda: supporting distributed programming in L++. In: EUROPEAN CONFERENCE ON OBJECT-ORIENTED PROGRAMMING, 7., 1993, Kaiserslautern. **Proceedings...** Berlin: Springer, 1993. p. 361–383. (Lecture Notes in Computer Science, v. 707).

CALLAGHAN, B.; PAWLOWSKI, B.; STAUBACH, P. **NFS version 3 protocol specification**: RFC 1831. London, 1995. 68 p.

CAROMEL, D.; KLAUSER, W.; VAYSSIÈRE, J. Towards seamless computing and metacomputing in Java. **Concurrency in Practice and Experience**, v. 10, n. 11–13, p. 1043–1061, sep / nov 1998. Disponível em: <<http://www-sop.inria.fr/~sloop/javall/index.ht>>. Acesso em: 20 fev. 2000.

FURMENTO, N.; ROUDIER, Y.; SIEGEL, G. **Parallélisme et distribution en C++**: une revue des langages existants. Valbonne, 1995. (RR 95-02). Disponível em: <<http://www-sop.inria.br/science/skd.gz>>. Acesso em: 29 fev. 2003.

JOEA, J. G.; JOHN, J. G. Importance of coffee in computer sciences. In: CONFERENCE ON COFFEE IMPORTANCE, 1., 2000, Java Island. **Proceedings...** Java Island: Java Island Press, 2003. p. 99–100.

LEE, K. fu. **AI Superpowers**. 1st. ed. New York: Houghton Mifflin Harcourt Publishing Company, 2018.

LENTINO, A. This chinese facial recognition start-up can identify a person in seconds. **CNBC**, 2019. Disponível em: <<https://www.cnbc.com/2019/05/16/this-chinese-facial-recognition-start-up-can-id-a-person-in-seconds.html>>. Acesso em: 23 mai. 2019.

LOUCKS DAVID SCHATSKY, T. D. J. State of ai in the enterprise, 2nd edition. **Deloitte Insights**, 2018. Disponível em: <[https://www2.deloitte.com/content/dam/insights/us/articles/4780\\_State-of-AI-in-the-enterprise/DI\\_State-of-AI-in-the-enterprise-2nd-ed.pdf](https://www2.deloitte.com/content/dam/insights/us/articles/4780_State-of-AI-in-the-enterprise/DI_State-of-AI-in-the-enterprise-2nd-ed.pdf)>. Acesso em: 23 mai. 2019.



MORGADO, M. L. C. **Reimplante dentário**. Trabalho de Conclusão de Curso (Especialização do curso) — Faculdade de Odontologia, Universidade Federal do Paraná, São Paulo, 2003.

NASCIMENTO, E. A. do. **Análise de curvas curvilíneas da trajetória da bola**. 1970. 36 f. Dissertação (Mestrado em Ciência do Futebol) — Cosmos University, Cidade do Cabo, 1971.

PATAGONIOS, J. **Um exemplo de TG**. 98 p. Trabalho de Conclusão de Curso (Graduação em Engenharia de Computação) — Instituto Teórico Aeroglifo, Santa Pindamonhangaba, 2001.

SBORNIAN, W. **Um exemplo de tese de doutorado**. 2004. 169 f. Tese (Doutorado em Aeronáutica) — Instituto de Alguma Coisa, Universidade Sei Lá de Onde, Santo Antônio da Patrulha, 2004. 1 CD-ROM.

ZYMERGEN: Early adopters combine bullish enthusiasm with strategic investments. 2019. Disponível em: <<https://www2.deloitte.com/insights/us/en/focus/cognitive-technologies/state-of-ai-and-intelligent-automation-in-business-survey.html>>. Acesso em: 23 may. 2019.

# Apêndice A - Tópicos de Dilema Linear

## A.1 Uma Primeira Seção para o Apêndice

A matriz de Dilema Linear  $M$  e o vetor de torques inerciais  $b$ , utilizados na simulação são calculados segundo a formulação abaixo:

$$M = \begin{bmatrix} M_{11} & M_{12} & M_{13} \\ M_{21} & M_{22} & M_{23} \\ M_{31} & M_{32} & M_{33} \end{bmatrix} \quad (\text{A.1})$$



FIGURA A.1 – Uma figura que está no apêndice

# Anexo A - Exemplo de um Primeiro Anexo

## A.1 Uma Seção do Primeiro Anexo

Algum texto na primeira seção do primeiro anexo.

## FOLHA DE REGISTRO DO DOCUMENTO

1. CLASSIFICAÇÃO/TIPO DM	2. DATA 25 de março de 2015	3. DOCUMENTO Nº DCTA/ITA/DM-018/2015	4. Nº DE PÁGINAS 34
5. TÍTULO E SUBTÍTULO: Uma Abordagem Sobre o Dilema do Cupim Frente ao Concreto Armado Utilizando Diferentes Composições Cimentísticas			
6. AUTORA(ES): <b>Maria das Graças Silva</b>			
7. INSTITUIÇÃO(ÕES)/ÓRGÃO(S) INTERNO(S)/DIVISÃO(ÕES): Instituto Tecnológico de Aeronáutica – ITA			
8. PALAVRAS-CHAVE SUGERIDAS PELA AUTORA: Cupim; Cimento; Estruturas			
9. PALAVRAS-CHAVE RESULTANTES DE INDEXAÇÃO: Cupim; Dilema; Construção			
10. APRESENTAÇÃO: <input checked="" type="checkbox"/> <b>Nacional</b> <input type="checkbox"/> <b>Internacional</b> ITA, São José dos Campos. Curso de Mestrado. Programa de Pós-Graduação em Engenharia Aeronáutica e Mecânica. Área de Sistemas Aeroespaciais e Mecatrônica. Orientador: Prof. Dr. Adalberto Santos Dupont. Coorientadora: Prof <sup>ª</sup> . Dr <sup>ª</sup> . Doralice Serra. Defesa em 05/03/2015. Publicada em 25/03/2015.			
11. RESUMO: <p>Aqui começa o resumo do referido trabalho. Não tenho a menor idéia do que colocar aqui. Sendo assim, vou inventar. Lá vai: Este trabalho apresenta uma metodologia de controle de posição das juntas passivas de um manipulador subatuado de uma maneira subótima. O termo subatuado se refere ao fato de que nem todas as juntas ou graus de liberdade do sistema são equipados com atuadores, o que ocorre na prática devido a falhas ou como resultado de projeto. As juntas passivas de manipuladores desse tipo são indiretamente controladas pelo movimento das juntas ativas usando as características de acoplamento da dinâmica de manipuladores. A utilização de redundância de atuação das juntas ativas permite a minimização de alguns critérios, como consumo de energia, por exemplo. Apesar da estrutura cinemática de manipuladores subatuados ser idêntica a do totalmente atuado, em geral suas características dinâmicas diferem devido a presença de juntas passivas. Assim, apresentamos a modelagem dinâmica de um manipulador subatuado e o conceito de índice de acoplamento. Este índice é utilizado na sequência de controle ótimo do manipulador. A hipótese de que o número de juntas ativas seja maior que o número de passivas (<math>n_a &gt; n_p</math>) permite o controle ótimo das juntas passivas, uma vez que na etapa de controle destas há mais entradas (torques nos atuadores das juntas ativas), que elementos a controlar (posição das juntas passivas).</p>			
12. GRAU DE SIGILO: <input checked="" type="checkbox"/> <b>OSTENSIVO</b> <input type="checkbox"/> <b>RESERVADO</b> <input type="checkbox"/> <b>SECRETO</b>			