

# João D. Ferreira

COMPUTER SCIENTIST · BIOINFORMATICIAN

R. das Flores, nº 7 – 1º dto, 2810-218 Laranjeiro, Portugal

☎ (+351) 96 323 50 69 | ✉ [jdferreira@fc.ul.pt](mailto:jdferreira@fc.ul.pt) | 🌐 <http://jdferreira.net>

ORCID 0000-0002-6900-3168 | 📷 [jdferreira](#) | 📄 [jdferreira](#)

## Summary

---

### Professional Experience

<b>2017–Present</b>	Integrated Member <i>at LASIGE (under the research line “Health and Biomedical Informatics”)</i>
<b>2015–Present</b>	Invited Assistant Professor <i>at Faculty of Sciences, University of Lisbon</i>
<b>2016</b>	Collaborator <i>at LASIGE (under the research line “Health and Biomedical Informatics”)</i>
<b>2012–2014</b>	Hired Researcher <i>by the EPIWORK project</i>

### Education

<b>2011–2016</b>	Ph.D. in Informatics, Specialization in Bioinformatics University of Lisbon
<b>2008–2010</b>	M.Sc. in Biochemistry Faculty of Sciences, University of Lisbon
<b>2005–2008</b>	B.Sc. in Biochemistry Faculty of Sciences, University of Lisbon

### Scientific Performance

<b>Publications</b>	7 peer-reviewed journal publications 10 conference and workshop publications total citations: 157 at the end of 2017 <i>h</i> -index: 7 (from Google Scholar)
<b>Knowledge resources</b>	2 knowledge-representation artefacts
<b>Projects</b>	participation in 3 national projects and 1 European project
<b>Academic Supervision</b>	1 Ph.D. in Statistics 1 M.Sc. in Bioinformatics & Computational Biology 2 M.Sc. in Informatics Engineering
<b>Conference organization</b>	ICBO 2015
<b>Reviewer</b>	3 international journals 2 international conferences 2 national conferences
<b>Honors and Awards</b>	2 scholarships (Ph.D. and M.Sc.)

### Other areas

My professional experience includes pedagogical performance, detailed in a dedicated section.

## A Scientific performance

### A.1 Scientific production

#### Selected list of publications

[2] Andre Lamurias, **João D. Ferreira**, and Francisco M. Couto. "Improving chemical entity recognition through h-index based semantic similarity." In: *Journal of Cheminformatics* 7.Suppl 1 Text mining for chemistry and the CHEMDNER track (Jan. 2015), S13. ISSN: 1758-2946. DOI: 10.1186/1758-2946-7-S1-S13

- **Impact Factor:** 4.550
- **Non-self citations:** 19

[4] Catia Pesquita, **João D. Ferreira**, Francisco M. Couto, and Mário J. Silva. "The epidemiology ontology: an ontology for the semantic annotation of epidemiological resources." In: *Journal of Biomedical Semantics* 5.1 (2014), p. 4. DOI: 10.1186/2041-1480-5-4

- **Impact Factor:** 2.262
- **Non-self citations:** 15

[5] **João D. Ferreira**, Janna Hastings, and Francisco M. Couto. "Exploiting disjointness axioms to improve semantic similarity measures." In: *Bioinformatics* 29.21 (2013), pp. 2781–2787. DOI: 10.1093/bioinformatics/btt491

- **Impact Factor:** 4.981
- **Non-self citations:** 20

[6] **João D. Ferreira**, Daniela Paolotti, Francisco M. Couto, and Mário J. Silva. "On the usefulness of ontologies in epidemiology research and practice." In: *Journal of epidemiology and community health* 0.0 (Nov. 2012). ISSN: 1470-2738. DOI: 10.1136/jech-2012-201142

- **Impact Factor:** 3.501
- **Non-self citations:** 10

[7] **João D. Ferreira** and Francisco M. Couto. "Semantic Similarity for Automatic Classification of Chemical Compounds". In: *PLoS Computational Biology* 6.9 (Sept. 2010). Ed. by John B. O. Mitchell, e1000937. ISSN: 1553-7358. DOI: 10.1371/journal.pcbi.1000937

- **Impact Factor:** 4.620
- **Non-self citations:** 43

#### Histograms of publications

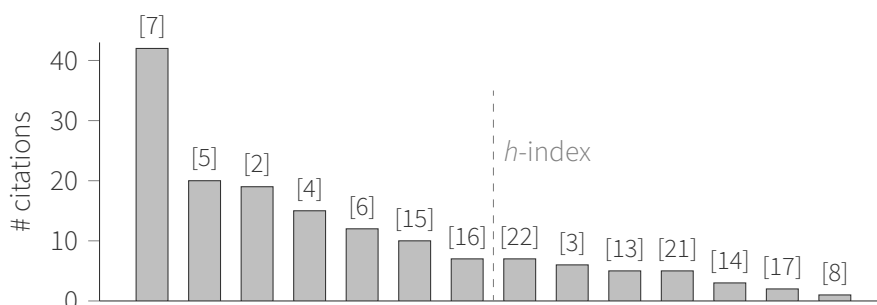


Figure 1: Histogram of citations for each of my cited papers, with papers sorted from most cited to least cited. Top labels specify the publication according to Section D. The vertical dashed line shows the associated *h*-index.

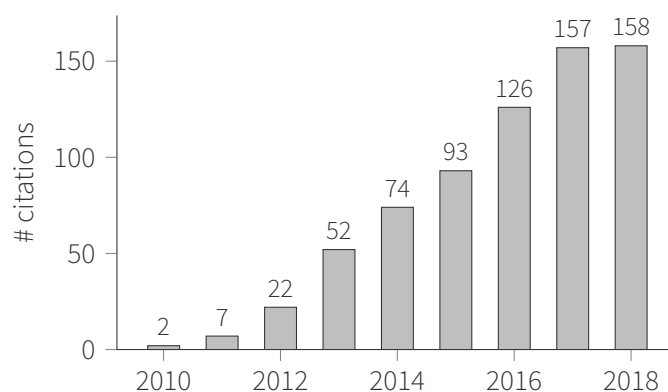


Figure 2: Histogram showing the cumulative number of citations gathered throughout the years.

## Selected list of citations

Stefan Schulz et al. “From concept representations to ontologies: a paradigm shift in health informatics?” In: *Healthcare informatics research* 19.4 (2013), pp. 235–242

- Mentions my paper [6] as **one of the twenty most influential papers** in the knowledge-representation, biomedical ontologies and electronic health issues, as listed by the members of the LinkedIn Working Group on Medical Concept Representation.

Janna Hastings et al. “The ChEBI reference database and ontology for biologically relevant chemistry: enhancements for 2013”. In: *Nucleic acids research* 41.D1 (2013), pp. D456–D463

- This paper is a major milestone in Chemoinformatics, describing ChEBI (an ontology of chemical compounds). It has itself 163 citations. Cites my paper [7] as an example of an **application that makes use of the semantic information** encoded in ChEBI for automatic classification.

Janna Hastings et al. “Structure-based classification and ontology in chemistry.” In: *J. Cheminformatics* 4 (2012), p. 8

- This paper expresses **the need and usefulness of systems that exploit the machine-readable information** provided by reference ontologies, and cites my paper [7] as an example.

Joerg Kurt Wegner et al. “Cheminformatics”. In: *Communications of the ACM* 55.11 (2012), pp. 65–75

- This paper advocates the use of **open-source solutions to deal with the vast amount of existing chemical information**, mentioning semantic similarity in ChEBI and citing my paper [7].

Robert Hoehndorf et al. “Thematic series on biomedical ontologies in JBMS: challenges and new directions.” In: *J. Biomedical Semantics* 5 (2014), p. 15

- This paper aims “to disseminate the latest developments in research on biomedical ontologies and provide a venue for publishing newly developed ontologies, updates to existing ontologies as well as methodological advances, and selected contributions from conferences and workshops”. It cites my paper [7] as an example of a semantic similarity application and my paper [5] as **the start of using more than taxonomies in semantic similarity** in the biomedical domain.

## Open Software

### OWLtoSQL

2015

URL: <http://github.com/jdferreira/owltosql>

- This software is one of the results of my PhD work. It is used to **convert an ontology written in OWL** (Web Ontology Language) **into a relational database**, which offers several advantages, the most promi-

ment of which is that it enables *random* access to the ontology constituents (classes, properties, axioms, etc.). It has been used by Bruno Inácio in his M.Sc work (see Section A.3).

### Multi-domain Ontology-based Semantic Similarity (MOSSy)

2015

URL: <http://github.com/jdferreira/mossy.bak>

- This software also resulted from my PhD work. It is responsible for calculating **semantic similarity between resources annotated with concepts from multiple ontologies**. The program is configurable to use any OWL ontology and can handle any type of annotated resource, both in single- and multi-domain contexts. It is extensible, since it allows the implementation of semantic similarity measures as Python classes.

## Knowledge representation artefacts

### Network of Epidemiology-Related Ontologies

2013

SEE MY PAPER [6]

- This work is a **collection of ontologies related to the epidemiology domain**. It contains 13 ontologies, from domains such as biochemistry, diseases, environment, transmission, vaccines, and geography, which provide the necessary concepts to annotate a digital resource of the epidemiological domain. The network was later included in the Epimarketplace, a repository of epidemic resources, in order to facilitate the annotation process.
- **Developed:** during my work in the EPIWORK project.

### Epidemiology Ontology

2014

URL: <https://code.google.com/p/epidemiology-ontology/>

- This ontology represents both **disease transmission methods and epidemiology parameters**, containing over 300 concepts. This ontology fills the previously existing gap in these areas, which were under-represented in biomedical ontologies. It supports the semantic annotation of resources.
- **Developed:** during my work in the EPIWORK project.

## A.2 Research projects

---

### Past projects

#### Semantic Mining with Linked Data (SMiLax)

2016-2019

URL: <http://lasige.di.fc.ul.pt/Projects/SMiLax>

- **Position:** Member of research team
- **Funding:** FCT (82 714€)
- **Grant:** PTDC/EEI-ESS/4633/2014
- **Description of project:** This project will improve the state of the art in semantic data mining by establishing novel methods and algorithms for the automated semantic annotation and enrichment of data, whose output will be explored by novel data mining approaches capable of capitalizing on the semantic web.
- **My contributions:** I am involved in research to improve current ontology-matching algorithms using semantic similarity instead of background knowledge.

#### Semantic Ontology Matching using External Resources (SOMER)

2012-2014

URL: <http://somer.rd.ciencias.ulisboa.pt/>

- **Position:** Member of research team
- **Funding:** FCT (84 000€)
- **Grant:** PTDC/EIA-EIA/119119/2010

- **Description of project:** The SOMER project aimed to develop ontology matching methods that exploit external knowledge resources through evidence and information content provided by unstructured text and annotation corpora. The results were applied to real world applications based on existing ontologies in the biomedical and geospatial areas.
- **My contributions:** I was involved in writing the proposal for this project. I was also responsible for exploring ways that allowed the use of semantic similarity measures in the process of ontology matching. I contributed to 1 peer-reviewed publication.

## Developing the Framework for an Epidemic Forecast Infrastructure (EPIWORK)

2009-2013

URL: <http://xldb.di.fc.ul.pt/wiki/EPIWork>

- **Position:** Hired researcher
- **Funding:** FP7 (5 000 000€)
- **Grant:** 231807
- **Description of project:** The EPIWORK project proposed a framework of tools and knowledge for the design of epidemic forecast infrastructures, including: development of the mathematical and computational methods needed to achieve prediction of disease spreading in complex social systems; development of large scale, data driven computational models aimed at epidemic scenario forecast; design and implementation of data-collection mechanisms, such as the collection of real-time disease incidence, through innovative web and ICT applications; and the implementation of a computational platform for epidemic research and data sharing.
- **My contributions:** The LASIGE team was responsible for the creation of the Epidemic Marketplace, an online platform for epidemic research and data sharing was the data-hub for the multiple research communities and countries involved in the project. I was hired by the LASIGE team as a consultant in semantic web, specifically to design the semantic metadata model that was used to annotate the resources of the marketplace, as well as to design the Network of Epidemiology-Related Ontologies, which contributed to the repository as a collection of concepts used to annotate the resources. I contributed to 5 peer-reviewed publications and 1 open-source ontology.

## Geographic Reasoning for Search Engines II (GREASE-II)

2010

- **Position:** M.Sc fellowship
- **Funding:** FCT (117 300 €)
- **Grant:** PTDC/EIA/73614/2006
- **Description of project:** This project aimed at researching information access methods to large collections of documents and objects having geographically rich text and meta-data. One of the main ideas of the project was that geospatial information can be mapped into ontology concepts.
- **My contributions:** One of the tasks of this project was to perform text-mining directly onto news reports in order to find geographical references in text; to assist in this task, a disambiguation module had to be created. I was assigned the creation of software to help disambiguate annotations manually, which would then later be used in a machine-learning step. I was also tasked with the alignment of Geo-Net-PT (a geographical ontology of the Portuguese territory) to Yahoo! GeoPlanet™ (an infrastructure for geo-referencing data on the Internet). I contributed to 1 technical report and to 1 open-source software.

## International Collaborations

### European Bioinformatics Institute (EBI)

2016-2017

- I collaborated with the team behind MetaboLights, a repository used to store details of experiments related to research in metabolism and information derived from the research. My role in this collaboration was to help **devise ways to measure annotation quality** in order to find resources under-annotated and to help authors submit properly annotated data. This collaboration resulted in 1 conference publication [8] and 1 peer-reviewed publication [1].

## University College London (UCL)

2013-2015

- I collaborated with Bernard D. Bono, a Principal Research Fellow in Health Informatics in UCL, in several fronts. The most recent collaboration is related to the **digital representation of functional tissue units** (small portions of an organ that perform a certain biological function) and the ways to compare them and publish them in a knowledge-base. This collaboration resulted in an oral presentation at the Virtual Physiological Human Conference 2014 [9].

## European Bioinformatics Institute (EBI)

2013

- I collaborated with J. Hastings, former group coordinator of the Chemoinformatics and Metabolism Team at EBI, in devising a semantic similarity measure that can effectively **capture the information provided by the disjointness axioms of an ontology**. This collaboration resulted in 1 peer-reviewed publication [5].

## Institute for Scientific Interchange (ISI) Foundation

2012

- I collaborated with Daniela Paolotti, the Project manager of the Italian web platform for Influenza-like Illness Surveillance. This collaboration resulted from my position at EPIWORK, and resulted in the publication of 1 peer-reviewed publication detailing **how to introduce semantic web technologies and ideas in the epidemiology domain** [6].

## A.3 Academic supervision

---

### Current supervision

#### PhD in Statistics

2017-2020

UNIVERSITY OF LISBON

- **Student:** Joana Stevens
- **Position:** Co-supervisor with Dr. Teresa Alpuim, University of Lisbon
- **Thesis:** “Application of machine-learning in vehicle insurance” (loosely translated from the official Portuguese title: “Aplicação de métodos de aprendizagem automática no desenvolvimento de um novo produto de seguro automóvel”)
- **Topic:** The work proposes the application of generalized linear models and machine learning algorithms in data from an insurance company regarding vehicle insurance. The final aim of this work is the establishment of customized insurance rates for every client, where the actual premium is calculated based on the values of a large number of variables not usually used in this area, thus creating an insurance policy that is more tailored to the client’s accident probability and more just to all clients in general.

### Past supervision

#### M.Sc in Informatics Engineering

2016-2017

UNIVERSITY OF LISBON

- **Student:** Isabela Mott
- **Position:** Co-supervisor with Dr Catia Pesquita, University of Lisbon
- **Thesis:** “AML-SSM: Semantic Similarity for Ontology Alignment”
- **Topic:** As part of my participation in the SMiLax project (see section A.2.1), I supervised the extension of the existing AgreementMaker Light (AML) platform, one of the most performant Ontology Alignment (OA) software pieces for biomedical ontologies. The student repurposed the idea of semantic similarity in order to further improve the existing OA algorithms in AML.

#### M.Sc in Informatics Engineering

2015-2016

UNIVERSITY OF LISBON

- **Student:** Bruno Inácio
- **Position:** Co-supervisor with Dr Francisco M. Couto, University of Lisbon
- **Thesis:** “How much is metadata worth?” (loosely translated from the official Portuguese title: “Quanto valem os metadados?”)
- **Topic:** The work focused (i) on the study of the quality of metadata as a way to sustain semantic integration and to facilitate data-sharing, based on the specificity of the ontology concepts used to annotate the resources and on the thoroughness of these annotations; and (ii) on the development of a platform that allows the assessment of metadata quality of the resources in a scientific data repository.

## Internship of a B.Sc in Biochemistry

2011-2012

UNIVERSITY OF LISBON

- **Student:** Hugo Ferreira
- **Position:** Co-supervisor with Dr Francisco M. Couto, University of Lisbon
- **Topic:** Creation of an ontology of epidemiology.
- **Topic:** The work consisted in using the textual descriptions in a Dictionary of Epidemiology to find relationships between the concepts and thus create an ontology of epidemiology.

## A.4 Scientific Dissemination

### Conference organisation

#### 6<sup>th</sup> International Conference on Biomedical Ontology (ICBO)

2015

WORKSHOP & TUTORIAL CHAIR, PROCEEDINGS CHAIR

- **Description:** The sixth International Conference on Biomedical Ontology (ICBO) was held in Lisbon in 2015 (<http://icbo2015.fc.ul.pt>). This prestigious and well-attended conference gathered multidisciplinary researchers at the University de Lisbon to present and discuss the latest research breakthroughs in exploring ontologies in a biomedical and clinical context.
- **Position:** Part of the organization; I was Workshop & tutorial chair and Session chair. I was also part of the Proceedings team, having compiled the Proceedings, and organized the submission to CEUR-WS.

### Reviewer

2018	Oxford Bioinformatics
2017	Journal of Biomedical Semantics
	Journal of Epidemiology & Community Health
2016	Journal of Biomedical Semantics
2015	Journal of Biomedical Semantics
	Journal of Epidemiology & Community Health
2013	Oxford Bioinformatics
2012	Oxford Bioinformatics

### Program Committee and other reviewer activities

2018	Conference on Practical Applications of Computational Biology & Bioinformatics
2017	Conference on Practical Applications of Computational Biology & Bioinformatics
2015	Bioinformatics Open Days
	Conference on Practical Applications of Computational Biology & Bioinformatics
	Portuguese Conference on Artificial Intelligence
2014	Conference on Practical Applications of Computational Biology & Bioinformatics
	International Conference on Data Integration & Life Sciences

## A.5 Grants, Honors and Awards

---

### Grants

- 2010–2015** PhD Grant by Fundação para a Ciência e Tecnologia (the Portuguese Science Funding Agency)
- 2009–2010** Pre-PhD grant by LASIGE

### Awards

- 2014** BPH Travel Award by the VPH Institute
- 2006** Best Student in 1<sup>st</sup> year by University of Lisbon



## B Pedagogical performance

### B.1 Teaching

#### Courses

Table 1: This table shows my contribution to courses from the Computer Science Department in Faculty of Sciences, University of Lisbon.

**Legend:** *PL*: Laboratory classes; *TP*: Practical lectures; *T*: Theoretical lectures

Course	2015–2016	2016–2017	2017–2018
IC	TP & PL	–	PL & TP
IPM	–	–	TP
ADSI	–	–	TP
ITW	TP & PL	TP & PL	TP & PL
ASW	PL	PL	–
PD	TP	TP	TP
Prog	TP	TP	–

#### IC (Computer Interaction)

2015–2018

- **Degree:** B.Sc. in Information Technologies (2<sup>nd</sup> year)
- **My classes:** TP & PL
- **Topics:**
  - Introduction to Human-Computer Interaction (HCI);
  - The foundations of HCI: Human and technological aspects;
  - The design process: user centred design, interaction design basics, guidelines for interaction design, and evaluation techniques;
  - Models and theories: cognitive models, task analysis, dialogue notations
- **My contributions:** Laboratory materials and project descriptions.

#### IPM (Human-Computer Interfaces)

2017–2018

- **Degree:** Informatics Engineering (2<sup>nd</sup> year)
- **My classes:** TP
- **Topics:**
  - The fundamental concepts of communication between humans and computers;
  - Technologies and styles of interaction;
  - Methods, principles and techniques for task analysis and design of interactive systems;
  - Usability evaluation;
  - The iterative cycle.
- **My contributions:** Laboratory materials and project descriptions.

#### ADSI (Information Systems Analysis and Design)

2017–2018

- **Degree:** Informatics Engineering (3<sup>rd</sup> year)
- **My classes:** TP
- **Topics:**
  - Data and functional analysis;
  - Analysis of human, organizational and environmental contexts in the context of information systems design;
  - “Understanding the problem” – emphasizing contact with the sources and analyzing a problem’s different perspectives;
  - “Developing a solution” – searching for an innovative, useful and well projected solution to the identified problem.

## ITW (Introduction to Web Technologies)

2015–2017

- **Degree:** B.Sc. in Information Technologies (1<sup>st</sup> year)
- **My classes:** TP & PL
- **Topics:**
  - The fundamental characteristics of the Web and associated technologies
  - The models and architectures that support the Web.
  - The main protocols (e.g. HTTP), specification and web programming languages (HTML, CSS, JavaScript, etc.)
  - Current platforms (e.g. W3.CSS, jQuery) that shape the Web.
- **My contributions:** Class materials used to teach HTML, CSS and JavaScript; design of project ideas based on the development of games as websites.

## ASW (Web Applications and Services)

2015–2017

- **Degree:** B.Sc. in Information Technologies (2<sup>nd</sup> year)
- **My classes:** Laboratory classes
- **Topics:**
  - Web application characteristics and features;
  - The development process of web applications;
  - Introduction to the main server-side web technologies: resource addressing, protocols and general architecture;
  - The various data transfer formats (XML, JSON, etc.) and related technologies;
  - Introduction to Web Services and Semantic Web.
- **My contributions:** Tutorial guides teaching how to deal with HTML forms and user input with Javascript.

## PD (Data Processing)

2015–2018

- **Degree:** B.Sc. in Biology (2<sup>nd</sup> year)
- **My classes:** TP
- **Topics:**
  - Introduction to data processing;
  - Introduction to Python (data types and data structures);
  - Introduction to Regular Expressions;
  - Introduction to Biomedical Web Services;
  - Database management systems.
- **My contributions:** I created all the practical lecture materials – a tutorial to guide students in processing large datasets, with emphasis on a collection of protein and metabolic information extracted from widely-known biomedical web services, including protein sequences and metabolic pathway data. The goal was to teach students how to process biology-related data with a programming language and an underlying database. Git repository in <https://github.com/jdferreira/data-processing-book>.

## Prog (Programming 1)

2015–2017

- **Degree:** Several B.Sc and M.Sc courses offered at Faculty of Sciences, University of Lisbon
- **My classes:** TP
- **Topics:**
  - Computation: computability and Turing machines;
  - Algorithms: exhaustive search, approximation search and bisection search;
  - Programming methods: attribution and verification, decision, iteration and recursion, abstraction and specification, cloning;
  - Programming languages: expressions and types, precedence and associativity, functions, scope, libraries and modules;
  - Data structures: sequences, tuples, lists and dictionaries;
  - Files;

- Software development: reading and writing, documentation, assertions and exceptions, test and debugging.
- **My contributions:** I contributed to the practice class materials by suggesting new exercises and modifications to existing ones.

## Pedagogical surveys

Lecturers, in Faculty of Sciences, University of Lisbon, are graded at the end of the semester by their students. These are my results. Results from the current year are not available yet.

Table 2: Pedagogical survey results. Students could either not answer each question or answer from 1 (strong disagreement) to 4 (strong agreement). Results for each question show the average for the students that answered the question. Results from the academic year of 2016–2017 are still unavailable.

### Legend:

Q1: Did the professor lecture with clarity?

Q2: Did the professor answer questions with clarity?

Q3: Was the professor available for outside-of-class contact & support?

Q4: Was there a good pedagogical relation between professor and students?

Q5: What is your global appreciation of the professor?

PL: Laboratory class;

TP: Practical lecture;

T: Theoretical lecture;

\*: Course still ongoing, results unavailable.

2015–2016	IC		Prog I	ASW	ITW	
Question	TP	PL	TP	PL	TP	PL
Q1	3.86	3.86	3.69	3.76	3.68	3.70
Q2	3.86	3.80	3.72	3.84	3.58	3.70
Q3	3.70	3.71	3.62	3.87	3.69	3.94
Q4	3.93	3.86	3.80	3.88	3.63	3.85
Q5	3.87	3.88	3.79	3.84	3.70	3.73

2016–2017	PD	Prog I	ASW	ITW	
Question	TP	TP	PL	TP	PL
Q1	3.42	3.53	3.53	3.72	3.57
Q2	3.50	3.61	3.59	3.71	3.75
Q3	3.39	3.57	3.36	3.74	3.66
Q4	3.45	3.55	3.50	3.74	3.71
Q5	3.44	3.55	3.55	3.79	3.75

## B.2 Jury and Examinations

### M.Sc in “Informatics Engineering”

2018

- **Student:** Pedro Davim Teixeira Mendes
- **Title:** “Development of custom solutions in Sharepoint” (loosely translated from the Portuguese official title “Desenvolvimento de soluções à medida em Sharepoint”)
- **Supervisor:** Dr Luis Antunes (Faculty of Sciences, University of Lisbon)

### M.Sc in “Informatics Engineering”

2017

- **Student:** Tiago Alexandre Fernandes de Noronha

- **Title:** “Mobile application development and support services” (loosely translated from the Portuguese official title “Desenvolvimento de aplicação móvel e serviços de suporte”)
- **Supervisor:** Dr João Balsa Silva (Faculty of Sciences, University of Lisbon)

### **M.Sc in “Information Systems and Computer Engineering”**

2017

- **Student:** Sebastião da Silva Freire
- **Title:** “E-Sports Ontology”
- **Supervisor:** Dr H Sofia Pinto (IST, University of Lisbon)
- I was formally invited by the supervisor; there has already been a discussion on the proposal of the M.Sc, which I examined in February 2017.

### **M.Sc in “Mathematics Applied to Economics and Business”**

2017

- **Student:** Catarina Nunes Valente
- **Title:** “Excel programming for Statistics: Linear Models and Extensions” (loosely translated from the Portuguese official “Programação em Excel para Estatística: Modelo Linear e Extensões”)
- **Supervisor:** Dr Teresa Alpuim (Faculty of Sciences, University of Lisbon)

### **M.Sc in “Bioinformatics and Computational Biology”**

2016

- **Student:** Samuel Viana
- **Title:** “Optimizing 16S Sequencing Analysis Pipelines”
- **Supervisors:** Daniel Faria (Instituto Gulbenkian de Ciência) and Catia Pesquita (Faculty of Sciences, University of Lisbon)

## **B.3 Teaching-related activities**

---

### **Invited participation in courses**

#### **Big Data**

2015-2017

- **Degree:** Ph.D. in Informatics
- I supervised the lectures and journal club for the course regarding the topics of machine learning in Big Data and the semantic web in Big Data.

#### **AW (Web Applications)**

2014–2015

- **Degree:** M.Sc. in Informatics Engineering
- I presented lectures on Semantic Web (SW) with the following topics:
  - The problem of ambiguity that SW tries to solve;
  - Rule-based inference;
  - RDF statements;
  - Several of the SW languages (RDF, OWL, SPARQL);
  - Objects vs. Classes vs. Instances
  - An introduction to several of the layers of the SW, including URIs, XML, RDF, Ontologies and Rules;
  - Real-world examples of SW in action: Semantic wikis, FOAF project, RDFa, hCalendar, Linked Data Project.

#### **Bioinformatics & Computational Modelling**

2013–2014

- **Degree:** Ph.D. Program in Biological Systems – Functional & Integrative Genomics
- I presented a practical lecture on Bioinformatics, specifically on the use of Python in biomedical data processing. This included:
  - Basic python datatypes and functions;
  - Introduction to BioPython, a package with access to several functions dedicated (i) to biological data processing and (ii) to widely-known biomedical web services;

- Exercises directed at learning the inners of BioPython, specifically to process protein sequences (using web services to access the SwissProt database and the BLAST algorithm)
- Introduction to the Gene Ontology and Semantic Similarity

## **Biomedical Ontologies**

2013–2014

- **Degree:** M.Sc. class available to several M.Sc. students at Faculty of Sciences, University of Lisbon
- I collaborated on the practical classes by designing a project and supervising a group of students in implementing the project's specifications. The projects were designed so that students would develop an intuition about ontology development, ontology matching, and semantic similarity. Two of these projects resulted in publications at national and international level.

## **Bioinformatics**

2011–2012

- **Degree:** B.Sc. class available to several B.Sc. students at Faculty of Sciences, University of Lisbon
- I was a teaching assistant on the practical classes by invitation of the course's head professor. I created a project specification based on ontology development and text-mining, and supervised the students in their implementations.

## C Ongoing Research

---

My current research efforts are focus on the application of semantic-web related technologies to multidisciplinary data, especially biomedical data, which usually exhibits high degree of multidisciplinary. I summarize my efforts in four topics.

### C.1 Multidisciplinary data-mining

---

Application of information-mining technologies to multidisciplinary data, including text- and data-mining of biomedical related information, with a particular emphasis on the use of semantic similarity measures to that effect. Data sources include health records, metabolic pathway descriptions, scientific publications, digital representations of biomedical entities, *etc.*

### C.2 Machine-learning in biomedical data

---

Application of machine-learning algorithms (in particular neural networks and deep-learning algorithms) to biomedical data. I have not yet started to research this topic, but I will in the very near future.

### C.3 Semantic similarity

---

Exploration of several ideas to improve semantic similarity measures in a multidisciplinary content. For the moment, this consists in a set of measures of relevance of concepts in an ontology.

### C.4 Data quality

---

Assessment of data quality by measuring some characteristics of the associated metadata, which including measuring the specificity and coverage of the metadata.

## D Full publication list

---

### D.1 Journal papers

---

- [1] **João D. Ferreira**, B Inácio, Reza M. Salek, and Francisco M. Couto. “Assessing Public Metabolomics Metadata, Towards Improving Quality.” In: *Journal of Integrative Bioinformatics* 14.4 (2017). doi: 10.1515/jib-2017-0054.
- [2] Andre Lamurias, **João D. Ferreira**, and Francisco M. Couto. “Improving chemical entity recognition through h-index based semantic similarity.” In: *Journal of Cheminformatics* 7.Suppl 1 Text mining for chemistry and the ChEMDNER track (Jan. 2015), S13. ISSN: 1758-2946. doi: 10.1186/1758-2946-7-S1-S13.
- [3] Andre Lamurias, **João D. Ferreira**, and Francisco M. Couto. “Identifying interactions between chemical entities in biomedical text”. In: *Journal of Interactive Bioinformatics (JIB)* 11.3 (2014), pp. 1–18. ISSN: 1613-4516. doi: 10.2390/biecoll-jib-2014-247.
- [4] Catia Pesquita, **João D. Ferreira**, Francisco M. Couto, and Mário J. Silva. “The epidemiology ontology: an ontology for the semantic annotation of epidemiological resources.” In: *Journal of Biomedical Semantics* 5.1 (2014), p. 4. doi: 10.1186/2041-1480-5-4.
- [5] **João D. Ferreira**, Janna Hastings, and Francisco M. Couto. “Exploiting disjointness axioms to improve semantic similarity measures.” In: *Bioinformatics* 29.21 (2013), pp. 2781–2787. doi: 10.1093/bioinformatics/btt491.
- [6] **João D. Ferreira**, Daniela Paolotti, Francisco M. Couto, and Mário J. Silva. “On the usefulness of ontologies in epidemiology research and practice.” In: *Journal of epidemiology and community health* 0.0 (Nov. 2012). ISSN: 1470-2738. doi: 10.1136/jech-2012-201142.
- [7] **João D. Ferreira** and Francisco M. Couto. “Semantic Similarity for Automatic Classification of Chemical Compounds”. In: *PLoS Computational Biology* 6.9 (Sept. 2010). Ed. by John B. O. Mitchell, e1000937. ISSN: 1553-7358. doi: 10.1371/journal.pcbi.1000937.

### D.2 Conferences (oral presentations & posters)

---

- [8] Bruno Inácio, **João D. Ferreira**, and Francisco M. Couto. “Metadata Analyser: Measuring Metadata Quality”. In: *11th International Conference on Practical Applications of Computational Biology & Bioinformatics*. Ed. by Florentino Fdez-Riverola, Mohd Saberi Mohamad, Miguel Rocha, Juan F. De Paz, and Tiago Pinto. Cham: Springer International Publishing, 2017, pp. 197–204. ISBN: 978-3-319-60816-7.
- [9] **João D. Ferreira**, Bernard de Bono, and Francisco M. Couto. “From data to knowledge: a tool for clustering multi-scale resources for physiology research”. In: *Virtual Physiological Human*. 2014.
- [10] Andre Lamurias, **João D. Ferreira**, and Francisco M. Couto. “Chemical named entity recognition: Improving recall using a comprehensive list of lexical features”. In: *8th International Conference on Practical Applications of Computational Biology & Bioinformatics*. Vol. 294. AISC. 2014, pp. 253–260. doi: 10.1007/978-3-319-07581-5\_30.
- [11] **João D. Ferreira**, Catia Pesquita, Mário J. Silva, and Francisco M. Couto. “Digital preservation of epidemic resources: coupling metadata and ontologies”. In: *International Conference on Preservation of Digital Objects*. Ed. by José Borbinha, Michael Nelson, and Steve Knight. 2013. URL: [http://xldb.di.fc.ul.pt/xldb/publications/Ferreira.etal:DigitalPreservationOf:2013\\_document.pdf](http://xldb.di.fc.ul.pt/xldb/publications/Ferreira.etal:DigitalPreservationOf:2013_document.pdf).
- [12] Catia Pesquita, **João D. Ferreira**, Francisco M. Couto, and Mário J. Silva. “Semi-automated annotation of epidemiological resources”. In: *Epiwork International Workshop “Digital Epidemiology”*. Epiwork International Workshop on Digital Epidemiology, 2013. URL: <http://webpages.fc.ul.pt/~fjcouto/files/abstract%20cpesquita-epiwork2013.pdf>.

- [13] David S. Batista, **João D. Ferreira**, Francisco M. Couto, and Mário J. Silva. “Toponym Disambiguation using Ontology-based Semantic Similarity”. In: *Computational Processing of the Portuguese Language*. Vol. 7243. 2012, pp. 179–185. ISBN: 9783642288845. DOI: 10.1007/978-3-642-28885-2\_20.
- [14] Francisco M. Couto, **João D. Ferreira**, João Zamite, Carlos Santos, Tiago Posse, Paulo Graça, Dulce Domingos, and Mário J. Silva. “The Epidemic Marketplace Platform: towards semantic characterization of epidemiological resources using biomedical ontologies”. In: *International Conference on Biomedical Ontologies*. 2012. URL: [http://ceur-ws.org/Vol-897/demo\\_1.pdf](http://ceur-ws.org/Vol-897/demo_1.pdf).
- [15] **João D. Ferreira**, Catia Pesquita, Francisco M. Couto, and Mário J. Silva. “Bringing epidemiology into the Semantic Web”. In: *International Conference on Biomedical Ontologies*. Vol. 897. 2012. URL: <http://ceur-ws.org/Vol-897/session1-paper02.pdf>.
- [16] **João D. Ferreira** and Francisco M. Couto. “Generic semantic relatedness measure for biomedical ontologies”. In: *International Conference on Biomedical Ontologies*. Vol. 833. 2011, pp. 117–123. URL: <http://ceur-ws.org/Vol-833/paper16.pdf>.
- [17] Bruno Tavares, Hugo Bastos, Daniel Faria, **João D. Ferreira**, Tiago Grego, Catia Pesquita, and Francisco M. Couto. “The Biomedical Ontology Applications (BOA) framework”. In: *International Conference on Biomedical Ontologies*. Vol. 833. 2011, pp. 300–301. URL: <http://ceur-ws.org/Vol-833/paper56.pdf>.

### D.3 Thesis (see section E)

---

- [18] **João D. Ferreira**. “Semantic Similarity Across Biomedical Ontologies”. PhD Thesis. Universidade de Lisboa, 2016.
- [19] **João D. Ferreira**. “Structural and semantic similarity metrics for chemical compound classification”. Master’s Dissertation. Universidade de Lisboa, 2010.

### D.4 Other

---

- [20] **João D. Ferreira**, Catia Pesquita, Francisco M. Couto, and Mário J. Silva. *Epiwork Deliverable 3.5: Epidemic Data Ontology*. Jan. 2012.
- [21] **João D. Ferreira**, David S. Batista, Francisco M. Couto, and Mário J. Silva. *The Geo-Net-PT / Yahoo! GeoPlanet TM concordance*. Tech. rep. October. Departamento de Informática – Faculdade de Ciências - Universidade de Lisboa, 2010. DOI: DOI:10455/6677.
- [22] Tiago Grego, **João D. Ferreira**, Catia Pesquita, Hugo Bastos, Diogo Vila Viçosa, João M. Freire, and Francisco M. Couto. *Chemical and Metabolic Pathway Semantic Similarity*. Tech. rep. Department of Informatics, Faculty of Sciences, University of Lisbon, 2010. DOI: DOI:10455/3335.



## E Education

---

### Ph.D. in Computer Science – Bioinformatics

Dec. 2010 – Jan. 2016

UNIVERSITY OF LISBON

- **Thesis:** Semantic Similarity Across Biomedical Ontologies
- **Supervisor:** Prof. Dr Francisco M. Couto
- **Grade:** Approved with Distinction and Honors
- **Abstract:** The need to compare complex entities is relevant in all the areas of science. In medicine, for example, comparing a clinical case to a database of previous cases can be extremely helpful when trying to diagnose a disease or deciding the most appropriate treatment for a patient.

Recent developments in knowledge representation, in particular the creation of the Web Ontology Language (OWL), have lead to a rise in the amount of knowledge that is being stored in *ontologies*, which represent, in machine-readable format, the known facts about reality. With the help of ontologies, statements like “**Influenza** is an **Infectious disease**” can be processed by computers, which, in turn, can be used to create new knowledge. In particular, *semantic similarity* has emerged to explore these ontologies as a way to compare entities annotated with the ontology concepts.

Semantic similarity has been extensively studied in the last decade, but some problems still persist. While there are algorithms to compare entities annotated with concepts from the same ontology, the possible ways to use *more than one ontology* are still in an early phase of study. For example, comparing a metabolic pathway using both the associated molecular functions and the metabolites converted in the pathway should, in principle, yield a higher precision than would be achieved with methodologies that rely on either one of the two domains independently. Comparing concepts from *different domains* and entities annotated with concepts from different domains is yet an unexplored area, but necessary to tackle multidisciplinary biomedical resources, e.g. to compare two clinical cases, the relationships between symptoms, diseases, blood screening results, etc. should provide a more insightful and precise value of similarity.

In this document, I explain the basic concepts needed to understand the problem of semantic similarity, how it is being solved, and how I propose to extend this notion so that it can be applied to more than one ontology and, more significantly, to more than one domain of knowledge.

### M.Sc. in Biochemistry

Sep. 2008 – June 2010

UNIVERSITY OF LISBON

- **Thesis:** Structural and semantic similarity metrics for chemical compound similarity
- **Supervisor:** Prof. Dr Francisco M. Couto
- **Grade:** 19 out of 20
- **Abstract:** Over the last few decades, there has been an increasing number of attempts at creating systems capable of comparing and classifying chemical compounds based on their structure and/or physicochemical properties. While the rate of success of these approaches has been increasing, particularly with the introduction of new and ever more sophisticated methods of machine learning, there is still room for improvement. One of the problems of these methods is that they fail to consider that similar molecules may have different roles in nature, or, to a lesser extend, that disparate molecules may have similar roles.

This thesis proposes the exploitation of the semantic properties of chemical compounds, as described in the CHEBI ontology, to create an efficient system able to automatically deal with the binary classification of chemical compounds. To that effect, I developed Chym (Chemical Hybrid Metric) as a tool that integrates structural and semantic information in a unique hybrid metric.

The work here presented shows substantial evidence supporting the effectiveness of Chym since it has outperformed all the models with which it was compared. Particularly, it achieved accuracy values of 90.9%, 87.7% and 84.2% when solving three classification problems which, previously, had only been solved with accuracy values of 81.5%, 80.6% and 82.8% respectively. Other results show that the tool is appropriate to use even if the problem at hand is not well represented in the CHEBI ontology. Thus, Chym shows that considering the semantic properties of a compound helps solving classification problems.

Therefore, Chym can be used in projects that require the classification and/or the comparison of chemical compounds, such as the study of the evolution of metabolic pathways, drug discovery or in preliminary toxicity analysis.

## **B.Sc. in Biochemistry**

*Sep. 2005 – June 2008*

UNIVERSITY OF LISBON

- **Discrimination of some grades:**

- Data Analysis and Processing in Biochemistry: 19
- Computational Biochemistry: 19
- Computational Simulation: 20

- **Final grade:** 18 out of 20