

David Soares Batista

PERSONAL DATA

PLACE AND DATE OF BIRTH: Lisbon, January 4th, 1981
PHONE: +49 (0)1523 850 98 89 (mobile)
EMAIL: dsbatista@inesc-id.pt
HOMEPAGE: <http://davidsbatista.github.io>

SUMMARY

I like to explore and extract knowledge from large volumes of data, be it structured or unstructured, numbers or text. I'm experienced in Machine Learning, Natural Language Processing and Big-Data.

In the past I've worked in projects, tackling problems with a strong text-mining and text-analysis components, involving tasks such as information extraction, classification, clustering and information retrieval.

In 2016, I successfully defended my Ph.D., proposing new methods to perform semantic relationship extraction from large collections of documents.

PROFESSIONAL EXPERIENCE

- | | |
|-----------------------|--|
| JAN. 2016 – PRESENT | Data Engineer @ HELLOFRESH AG (http://www.hellofresh.de) <ul style="list-style-type: none">• Text analysis of customer NPS comments with NLP and ML.• Data Warehousing and processes automation for data analysis.• Design, build and maintain ETLs and infrastructure for batch processing.• Technologies: Python NLTK, scikit-learn, Spark, Hive, Impala, Airflow. |
| JUN. 2011 – APR. 2014 | Researcher and Developer @ INESC-ID (http://www.inesc-id.pt)
Project: REACTION - Computational Journalism <ul style="list-style-type: none">• Named Entity Recognition and Disambiguation over news articles.• Semantic Relationship Extraction between named-entities.• Topic Modeling (i.e., LDA) techniques applied to on-line news.• Analysis of large graphs connecting entities and topics. |
| OCT. 2009 – OCT. 2011 | Researcher and Developer @ LASIGE (http://www.lasige.di.fc.ul.pt)
Project: GREASE - Geographic Reasoning for Search Engines <ul style="list-style-type: none">• Language identification on a web-crawl using n-grams models.• Disambiguation of geographic ambiguous names in Portuguese text.• Alignment between geo-ontologies by linking geographical references. |
| OCT. 2007 – JUL. 2008 | Software Developer @ NOKIA SIEMENS NETWORKS <ul style="list-style-type: none">• Development of data collection modules for a GSM monitoring system.• Technologies: Java, CORBA Architecture, Oracle RDBMS |

COMPUTER SKILLS

Programming:	Python, Java, Shell Scripting
Machine Learning:	scikit-learn, gensim, MinorThird
Natural Language Processing:	Python NLTK, Stanford CoreNLP, LingPipe
Distributed Computing:	Hadoop, Spark, Hive, Impala
Information Retrieval:	Apache Lucene/Solr
Databases/NoSQL:	MySQL, PostgreSQL, REDIS

EDUCATION

- 2011-2015 | Doctor of Philosophy (Ph.D.) **Instituto Superior Técnico**, University Lisbon.
“**Large-Scale Semantic Relationship Extraction**”
To achieve scalable relationship extraction, I proposed using an on-line classifier, based on the idea of nearest neighbor classification, and leveraging min-hash and locality sensitive hashing for efficient similarity search. To obtain training data, for the classifier, I proposed a bootstrapping technique relying on distributional word representations.
- 2007-2009 | Master’s Degree (M.Sc.) - **Faculty of Sciences**, University of Lisbon.
“**Geographic Text Mining**”
I developed an information extraction system based on Conditional Random Fields to generate geographic summaries. The summary lists the geographic entities found in a document and mapped into geographic concepts in a geographic ontology. The system was applied to a crawl of the Portuguese Web (25GB raw text) using a Hadoop cluster.
- 2003-2007 | Bachelor of Science (B.Sc.), Informatics Engineering,
Faculty of Sciences, University Lisbon.
- 2005-2006, **Karlsruhe Universität (TH)**, Germany
- Erasmus exchange student for two semesters
- 2004 NOVEMBER - 2005 JULY - IT support (part-time) @ University of Lisbon
- Troubleshooting network connections and services, preventative maintenance, helping and educating new users.
- 2003 NOVEMBER - 2004 MARCH - SysAdmin (part-time) @ University of Lisbon
- Administration and configuration of networking software and services (e.g., samba, backups, quotas, mail, web, crontab, IPTables).

LANGUAGES

PORTUGUESE:	<i>Native Speaker</i>		
ENGLISH:	Spoken: <i>Fluent</i>	Written: <i>Fluent</i>	
GERMAN:	Spoken: <i>Good</i>	Written: <i>Fair</i>	Certification: <i>Goethe-Zertifikat B.1</i>

SELECTED PUBLICATIONS

Semi-Supervised Bootstrapping of Relationship Extractors with Distributional Semantics. David S Batista, Bruno Martins, and Mário J Silva. In *Empirical Methods in Natural Language Processing-EMNLP’15*. - **Honorable Mention for Best Short Paper**

A Minwise Hashing Method for Addressing Relationship Extraction from Text. David S Batista, Rui Silva, Bruno Martins, and Mário J Silva. In *Web Information Systems Engineering-WISE’13*.

Exploring DBpedia and Wikipedia for Portuguese Semantic Relationship Extraction David Soares Batista, David Forte, Rui Silva, Bruno Martins, and Mário J. Silva. *Linguamática*, 5(1), 2013.

Toponym Disambiguation using Ontology-based Semantic Similarity David S Batista, João D Ferreira, Francisco M Couto, and Mário J Silva. In *Computational Processing of the Portuguese Language 2012*.

A Statistical Study of the WPT05 Crawl of the Portuguese Web David Batista and Mário J. Silva. In *In FALA 2010 VI Jornadas en Tecnología del Habla and II Iberian SLTech Workshop* Universidade de Vigo, 2010.

Geographic Signatures for Semantic Retrieval David S Batista, Mário J Silva, Francisco M Couto, and Bibek Behera. In *Proceedings of the 6th Workshop on Geographic Information Retrieval* ACM, 2010.