

HEIMDALLR

THE GUARDIAN API FOR ACADEMIC INTEGRITY

BY TOMAS SANCHEZ

FOR UNIVERSIDAD TECNOLÓGICA NACIONAL, FACULTAD REGIONAL BUENOS AIRES

NATURAL LANGUAGE PROCESSING, K3551

AS TAUGHT BY ENG. HERNAN BORRE

NOVEMBER 10TH, 2023

TABLE OF CONTENTS

TABLE OF CONTENTS	2
ABSTRACT	3
INTRODUCTION	4
FEATURES	4
METHOD	5
DOMAIN MODELING	5
DRIVEN ADAPTERS	5
EVENT-DRIVEN ORCHESTRATION	6
PLAGIARISM DETECTION AS A SERVICE	8
RESULTS	9
PERFORMANCE ISSUES	9
TRAINING DATA	ERROR! BOOKMARK NOT DEFINED.
METRICS ANALYSIS	11
ABOUT THIS DOCUMENT	12
LICENSE	13
AUTHOR	13
REFERENCES	14
APENDIX A	15

ABSTRACT

IN THIS RESEARCH, A PRODUCTION-READY REST API WAS DEVELOPED TO PRESENT PLAGIARISM DETECTION AS A SERVICE. FOR THIS, THE LATEST, MOST MODERN, ASYNCHRONOUS COMPATIBLE PYTHON FRAMEWORK WAS USED IN CONJUNCTION WITH PRE-TRAINED NATURAL LANGUAGE PROCESSING (NLP) MODELS FOR SPANISH TEXTS. WRITINGS ARE COMPARED BY SENTENCE AGAINST PREVIOUSLY STORED SENTENCES FROM DIFFERENT WORKS. PLAGIARISM IS DETERMINED BY COSINE SIMILARITY, WITH A CONFIGURABLE THRESHOLD.

KEYWORDS: NLP, REST API, PYTHON, COSINE SIMILARITY, ASYNCHRONOUS PROCESSING, COMMAND QUERY RESPONSIBILITY SEGREGATION (CQRS), EVENT DRIVEN, PLAGIARISM DETECTION.

INTRODUCTION

HEIMDALLR IS A POWERFUL AND INTELLIGENT PLAGIARISM DETECTION SYSTEM DESIGNED TO UPHOLD ACADEMIC INTEGRITY AND ENSURE THE AUTHENTICITY OF WRITTEN WORKS.

INSPIRED BY THE MYTHOLOGICAL FIGURE HEIMDALLR, THE ALL-SEEING, KNOWN FOR HIS UNWAVERING VIGILANCE. ITS ADVANCED DETECTION CAPABILITIES AND USER-FRIENDLY INTERFACE RESULT IN THE ULTIMATE SOLUTION FOR EDUCATIONAL INSTITUTIONS, RESEARCHERS, AND STUDENTS TO MAINTAIN THE HIGHEST STANDARDS OF ACADEMIC HONESTY AND ORIGINALITY.

THIS PLAGIARISM DETECTOR APPLICATION WAS DEVELOPED AS AN ASSIGNMENT FOR NATURAL LANGUAGE PROCESSING COURSE FROM UNIVERSIDAD TECNOLÓGICA NACIONAL, FACULTAD REGIONAL BUENOS AIRES DURING THE SECOND SEMESTER OF 2023.

FEATURES

- **COMPREHENSIVE PLAGIARISM DETECTION:** HEIMDALLR SCANS THROUGH A VAST ARRAY OF SOURCES, INCLUDING WEB CONTENT, BOOKS, PAPERS, PHOTOGRAPHS, AND PREVIOUS STUDENT WORKS, LEAVING NO ROOM FOR UNDETECTED PLAGIARISM.
- **PRECISE LOCATION IDENTIFICATION:** IT DOESN'T JUST FLAG POTENTIAL INSTANCES OF PLAGIARISM; HEIMDALLR PROVIDES PINPOINT ACCURACY BY HIGHLIGHTING THE EXACT LINES OR WORDS IN QUESTION, ENABLING USERS TO SWIFTLY ADDRESS THE ISSUE.
- **SOURCE ATTRIBUTION:** THE SYSTEM NOT ONLY DETECTS PLAGIARISM BUT ALSO IDENTIFIES THE SOURCE THAT IS BEING PLAGIARIZED.

METHOD

HEIMDALLR WAS BUILT FOLLOWING AN EVENT-ORIENTED ARCHITECTURE WITH SUPPORT FOR DOMAIN MODELING. PRESENTING A CLEAN ARCHITECTURE, ALLOWING CHANGES IN EACH LAYER HAS MINIMAL EFFECT ON OTHERS. THE ARRANGEMENT OF THIS DIVISION IS EASILY IDENTIFIED IN THE PACKAGE NAMING:

1. DOMAIN – THE CORE BUSINESS OBJECTS.
2. ADAPTERS – I/O ABSTRACTIONS.
3. SERVICE – EVENT HANDLING ORCHESTRATION
4. ENTRY POINT –APPLICATION DRIVERS.

DOMAIN MODELING

THE PROBLEM THAT THE APPLICATION TRIES TO SOLVE IS PLAGIARISM IN STUDENTS' ASSIGNMENTS, BEING THE MAIN ENTITY MODELED IN AN ASSIGNMENT DATA CLASS. THIS RECORD REPRESENTS PERSISTENT WRITINGS, SAVED IN A CONVENIENT FORMAT THAT FACILITATES ITS IDENTIFICATION, AND COMPARISON.

```
class Assignment (BaseDocument):
    title: str
    author: str
    content: list[str]
    date: datetime.date | None
    similarities: list[AssignmentVerification] | None
```

NOTICE THE USE OF TYPE ANNOTATIONS OFFER BETTER LEGIBILITY. THIS MODEL WILL BE PERSISTED IN A DATABASE. THE `BaseDocument` SUPERCLASS CONTAINS AN AUTO-GENERATED IDENTIFIER USING `UUID4` FORMAT.

DRIVEN ADAPTERS

INWARD-FACING ADAPTERS PRESENT A LAYER OF ABSTRACTION FOR READING, AND WRITING OPERATIONS PERFORMED OVER A DATABASE AND FILES.

AN `AssignmentReader` INTERFACE IS RESPONSIBLE FOR MAPPING TEXT INTO A DOMAIN MODEL. WITH ONLY ONE NAMED METHOD, GUARANTEES, NO MATTER WHAT IMPLEMENTATION IS DONE, THAT A PERSISTENT ENTITY CAN BE OBTAINED.

```
class AssignmentReader (abc.ABC):
    @abc.abstractmethod
    def read(file: UploadFile) -> Assignment:
        raise NotImplementedError
```

THE CURRENT INTERFACE IMPLEMENTATION UTILIZES `spaCy` FOR SENTENCE RECOGNITION: AN ASSIGNMENT CONTENT CONSISTS OF A LIST OF SENTENCES. ALSO, THIS NLP LIBRARY HELPS WITH AUTHOR RECOGNITION USING TOKENIZATION AND TOKEN TYPE FOR EACH SENTENCE. FOR EACH DIFFERENT CONTENT TYPE SUPPORTED, A DEPENDENCY IS NEEDED. `PYMuPDF`, `PYTEXTTRACT` AND `DOCX2TXT` ARE USED FOR TEXT PROCESSING.

FOR CONVENIENCE AND FAMILIARITY, THE DATABASE CHOSEN WAS MONGODB. WITH MOTOR ENGINE, NON-BLOCKING I/O OPERATIONS ARE PERFORMED. NOT DELVING INTO IMPLEMENTATION DETAILS, A REPOSITORY PATTERN IS USED FOR BASIC CRUD OPERATIONS.

EVENT-DRIVEN ORCHESTRATION

AN ORCHESTRATION LAYER IS NEEDED FOR EXECUTING OUR SERVICE VALUE: CHECKING PLAGIARISM IN WRITINGS. SERVICES RECEIVE AN INPUT AND PRODUCE A VALUE. THESE INPUTS ARE ORDERS ORIGINATED FROM ANOTHER SOURCE, IT COULD BE ANOTHER APPLICATION SERVICE, OR AN EXTERNAL SOURCE. REPRESENTED BY DATA CLASSES, **COMMANDS**. COMMANDS EXPECT THE SYSTEM TO PRODUCE AN EFFECT, SOMETHING MUST HAPPEN. ONCE THE SYSTEM REACTS TO A COMMAND, A NOTIFICATION IS PRODUCED, AN **EVENT**. THIS IS VISIBLE IN THE SERVICE INTERFACE.

```
class AssignmentVerifier (abc.ABC):
    @abc.abstractmethod
    async def veriify(command: VerifyAssignment) -> AssignmentVerified:
        raise NotImplementedError
```

NOTICE THE NAMING OF THE DATA CLASSES. COMMANDS ARE NAMED AFTER IN IMPERATIVE MOOD, CLEARLY REPRESENTING AN ORDER. MEANWHILE, EVENTS ARE NAMED IN PAST TENSE, SOMETHING THAT HAS HAPPENED.

AS AN ORCHESTRATOR, THE IMPLEMENTATION MAKES USE OF THE AFOREMENTIONED ABSTRACTIONS FOR FILE MAPPING, AND DATABASE OPERATIONS.

THE CODE BLOCK BELOW IS NOT THE FINAL IMPLEMENTATION BUT A HIGH-LEVEL EXAMPLE OF HOW A SERVICE CAN BE DEVELOPED DEPENDENT ON ADAPTERS. IT USES *spaCy*'S IMPLEMENTATION OF COSINE SIMILARITY FOR EACH SENTENCE.

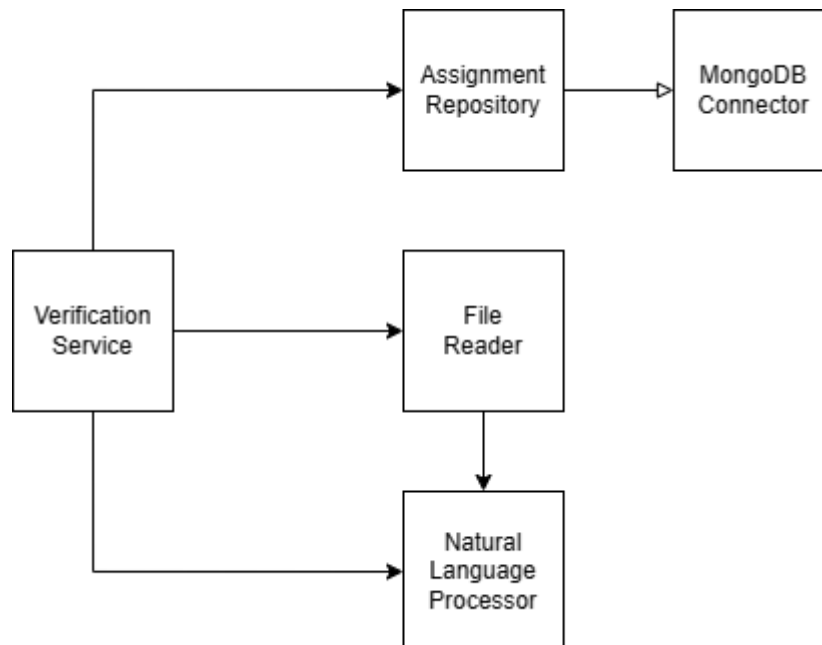


FIGURE 1 – ASSIGNMENT VERIFICATION SERVICE

```
async def verify(self, command: VerifyAssignment) -> AssignmentVerified:

    # use the AssignmentReader adapter to obtain a domain object
    entry = self.reader.read(file=command.file)
    entry.id = command.id

    # use the Repository/DB adapter to obtain all previous writings
    assignments: list[Assignment] = await self.repository.find_all()
    comparisons: list[AssignmentCompared] = []

    # compare entries with each stored assignments
    for assignment in assignments:
        # compare assignments compares each sentence
        result = self.compare_assignments(assignment, entry)
        if result.similarities:
            comparisons.append(result)

    # map the comparison results to the AssignmentVerification model
    verifications = [AssignmentVerification(**comparison.model_dump())
                     for comparison in comparisons]
    entry.similarities = verifications

    # use the DB adapters to persist the operation
    persisted = await self.repository.save(entry)

    # notifies the system an event has occurred
    return AssignmentVerified(
        id=persisted.id,
        title=persisted.title,
        author=persisted.author,
        similarities=comparisons,
    )

def compare_sentence(
    self,
    sentence: str,
    entry_sentence: str) -> SentenceCompared:
    persisted_sentence_doc = self.nlp(sentence)
    entry_sentence_doc = self.nlp(entry_sentence)
    # cosine similarity percentage
    similarity = persisted_sentence_doc.similarity(entry_sentence_doc)

    return SentenceCompared(
        present=sentence,
        compared=entry_sentence,
        plagiarism=similarity)
```

PLAGIARISM DETECTION AS A SERVICE

THE ADDED VALUE OF HEIMDALLR IS PROVIDING NLP EVALUATION THROUGH AN API. FOR THAT REASON, FASTAPI MICROFRAMEWORK WAS ADOPTED FOR DEVELOPING OUTWARD-FACING ADAPTERS, APPLICATION'S DRIVERS.

FASTAPI OFFERS HIGH PERFORMANCE, HIGH DEVELOPMENT SPEED, EASY AND INTUITIVE LINES OF CODE, AND IS STANDARDS-BASED. IT FOLLOWS OPENAPI ¹ AND JSON SCHEMA ².

IN THE AUTO-GENERATED DOCUMENTATION PAGE, WE CAN DISTINGUISH DIFFERENT PATH GROUPS.

1. **MONITOR** — USED FOR CHECKING APPLICATION AVAILABILITY.
2. **QUERIES** — REQUEST FOR INFORMATION ABOUT A SYSTEM. DOES NOT MAKE ANY CHANGES.
3. **COMMANDS** - REQUEST FOR A SPECIFIC ACTION TO OCCUR IN THE SYSTEM.

THE UI PROVIDES EXAMPLES OF BOTH REQUEST AND RESPONSE, AND IT ALSO DESCRIBES THE USED SCHEMAS. THESE ENTRY POINTS MAY REUTILIZE BOTH SERVICE AND ADAPTER ABSTRACTIONS TO PERFORM REQUESTED OPERATIONS.

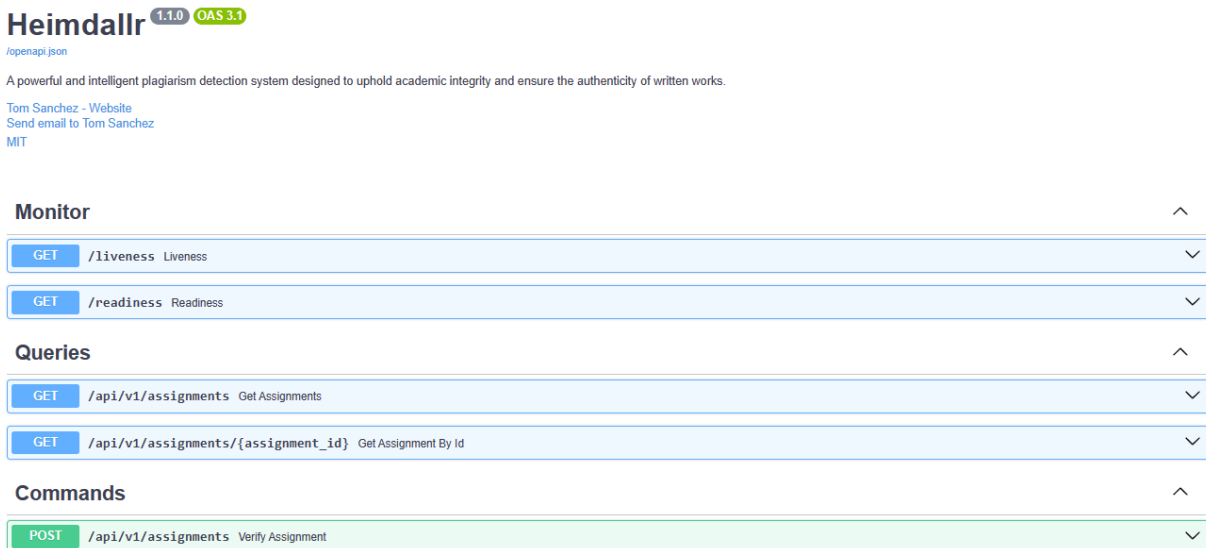


FIGURE 2 OPENAPI UI

¹ THE OPENAPI SPECIFICATION (OAS) DEFINES A STANDARD, PROGRAMMING LANGUAGE-AGNOSTIC INTERFACE DESCRIPTION FOR HTTP APIs.

² JSON SCHEMA IS THE VOCABULARY THAT ENABLES JSON DATA CONSISTENCY, VALIDITY, AND INTEROPERABILITY AT SCALE.

RESULTS

PERFORMANCE ISSUES

THE IDEA WAS TO VERIFY EACH SENTENCE OF A NEW ENTRY AGAINST ALL OTHER SENTENCES OF ALL STORED ASSIGNMENT. AS RESPONSE TIME STARTED TO INCREASE, SOLUTIONS TO NEEDED TO START POPPING UP.

AT FIRST, IT WAS TRIED TO BE FIXED AS AN ASYNCHRONOUS ENDPOINT: WHEN A FILE IS SUBMITTED, A UNIQUE UUID IS RETURNED. WITH THIS UUID, THE ASSIGNMENT CAN BE LATER RETRIEVED TO SEE ITS PLAGIARISM RESULTS. THIS SOLVED CLIENT TIMING OUT WAITING FOR SERVER RESPONSE ISSUE. HOWEVER, PROCESSING TIME WAS STILL TOO EXPENSIVE.

ANOTHER IDEA WAS TO ONLY COMPARE LONG SENTENCES. IT WAS RESEARCHED THAT ON AVERAGE, SPANISH HAS 5.6 AVERAGE LETTERS PER WORD³, AND ARBITRARILY DECIDED THAT USING 3 WORDS WITH LESS CHARACTERS THAN THE AVERAGE ISN'T BEING PLAGIARIZED. THIS HAD LITTLE TO ZERO IMPACT. COMPARING ONE ASSIGNMENT WITH ONLY 14 OTHERS WILL TAKE UP TO 15 MINUTES.

FINALLY, A PRIOR DISCARD IDEA CAME UP. COMPARING WHOLE DOCUMENTS, NOT SENTENCE BY SENTENCE. FOR THAT, A THRESHOLD NEEDED TO BE DEFINED. KNOWN ASSIGNMENTS WITH THE SAME REQUIREMENTS WERE MANUALLY COMPARED. USING "WIKINOMICS" ASSIGNMENT. WHEN EVALUATING MELANIE BLEJTER'S WORK WITH A THRESHOLD OF 0.99.

AUTHOR	SAME ASSIGNMENT	SIMILARITY	SECONDS
JOURDAN MARTIN	YES	0.9929278646150054	71.27109
MARCELO BROSSI	YES	0.9935113489761744	90.099331
IVO URSINO	NO	0.9901813848392802	89.807355

IT WAS SEEN THAT UNRELATED ASSIGNMENTS WERE ALSO HAVING SIMILARITY HIGHER THAN 0.99. HAVING UPDATED THRESHOLD TO 0.991 TESTING AGAIN WITH SANTIAGO PERALTA'S TAKE ON "WIKINOMICS".

AUTHOR	SAME ASSIGNMENT	SIMILARITY	SECONDS
JOURDAN MARTIN	YES	0.9922299777151997	71.809523
MARCELO BROSSI	YES	0.991994420219641	89.406367
MELANIE BLEJTER	YES	0.9956091815513386	90.936489

WITHOUT CHANGING ANYTHING, FOR MIGUEL ARTEAGA OROZCO'S WRITING.

AUTHOR	SAME ASSIGNMENT	SIMILARITY	SECONDS
SANTIAGO PERALTA	YES	0.9941490039858188	104.084537
MARCELO BROSSI	YES	0.991522754728401	114.398007
MELANIE BLEJTER	YES	0.9952789647155554	111.424194

THERE WASN'T ENOUGH SIMILARITY WITH JOURDAN MARTIN'S ASSIGNMENT. MANUALLY CHECKING IT, IT HAS 0.9896672681699993.

A THRESHOLD OF 0.991 WAS CONSIDERED ACCEPTABLE.

³ SEE ALSO: [WHY DOES CONTENT EXPAND WHEN TRANSLATED?](#) — ON [INTER CONTACT TRANSLATIONS'](#) BLOG.

TOPIC CLASSIFICATION

ANOTHER CHALLENGE WAS TO CATEGORIZE ASSIGNMENTS. HAVING MANUALLY LABELED 100 ASSIGNMENTS WITH DIFFERENT TOPICS, A *TRAIN TEST SPLIT* WITH A SIZE OF 0.33 AND RANDOM STATE OF 42 WAS DONE FOR TEXT CLASSIFICATION.

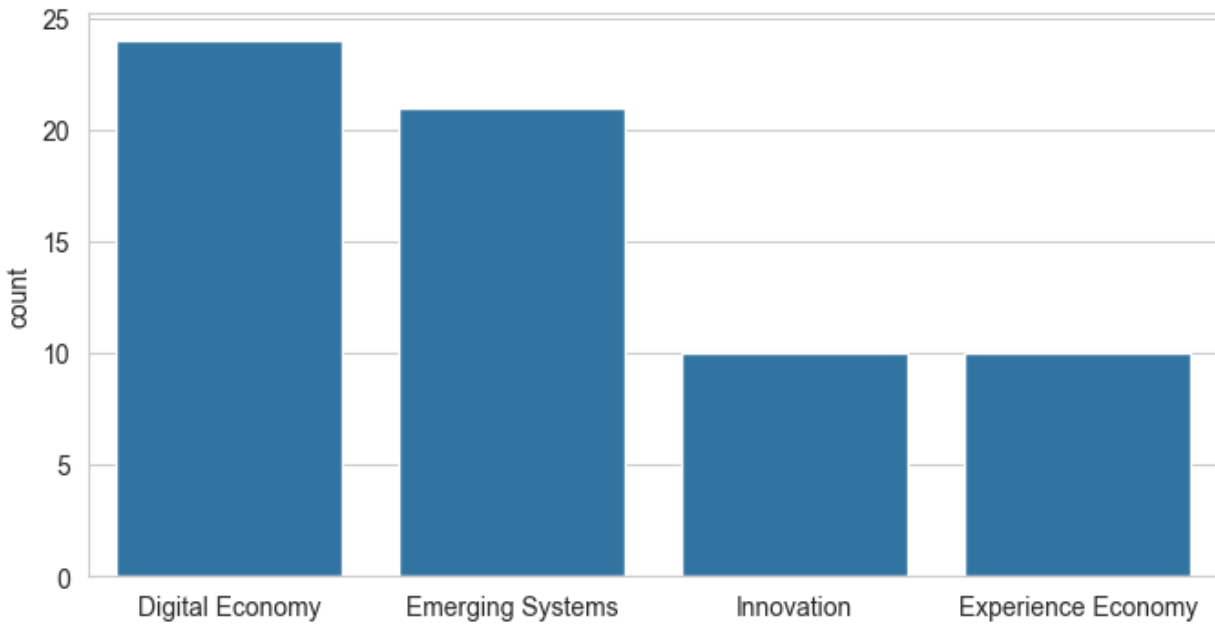


FIGURE 3 TEST SPLIT BY TOPIC

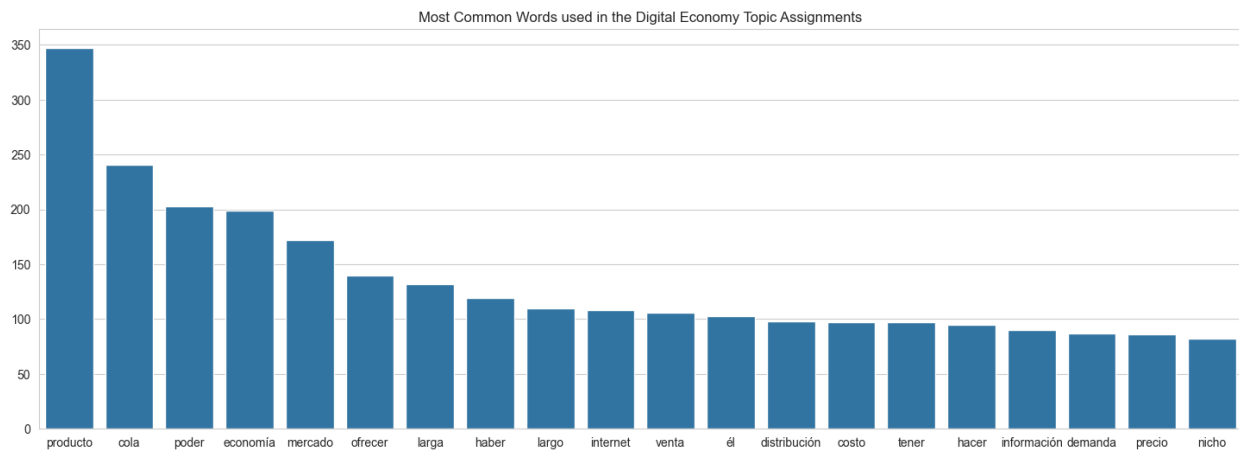


FIGURE 4 MOST COMMON WORDS IN DIGITAL ECONOMY

MORE GRAPHS CAN BE SEEN ON THE INCLUDED APPENDIX.

IT IS SEEN THAT “PRODUCTO” AND “COLA” ARE IN THE TOP WORDS FOR “DIGITAL ECONOMY”. IT IS OBVIOUS THAT “THE LONG TAIL” ASSIGNMENT WAS LABELED AS “DIGITAL ECONOMY”.

THE TOP WORDS FOR “INNOVATION” ARE “ADOPCION” AND “PRODUCTO”.

METRICS ANALYSIS

ACCURACY: 0.96969696969697

TOP 10 FEATURES USED TO PREDICT.

CLASS 1 BEST:

(-0.02910298582999386, 'EJEMPLO')

(-0.019688603502727383, 'E.')

(-0.01691601764291971, 'PODER')

(-0.01644621463970039, 'INDIVIDUO')

(-0.015634058315695205, 'SISTEMA')

(-0.014470155004150373, 'IDEA')

(-0.014373718553172994, 'CLIENTE')

(-0.014351650689569237, 'FORMA')

(-0.013436612415077655, 'PREGUNTAR')

(-0.01332442548778907, 'DESARROLLAR')

CLASS 2 BEST:

(0.04238784680702782, 'ECONOMÍA')

(0.029182893807513717, 'DIGITAL')

(0.027224404983200518, 'FUENTE')

(0.02054267956495503, 'ENERGÍA')

(0.019831701811292908, 'INTERNET')

(0.017813780608358233, 'COLA')

(0.015380094727189587, 'LARGO')

(0.015285828149245853, 'RED')

(0.01245561886065532, 'BIEN')

(0.012392980115877975, 'CUALQUIERA')

	PRECISION	RECALL	F-1 SCORE	SUPPORT
EMERGING SYSTEMS	1.00	0.92	0.96	12
EXPERIENCE ECONOMY	1.00	1.00	1.00	7
DIGITAL ECONOMY	1.00	1.00	1.00	5
INNOVATION	0.90	1.00	0.95	9
ACCURACY			0.97	33
MACRO AVG	0.97	0.98	0.98	33
WEIGHT AVG	0.97	0.97	0.97	33

TOPIC SIMILARITIES

USING 98 MANUALLY LABELLED ASSIGNMENTS.

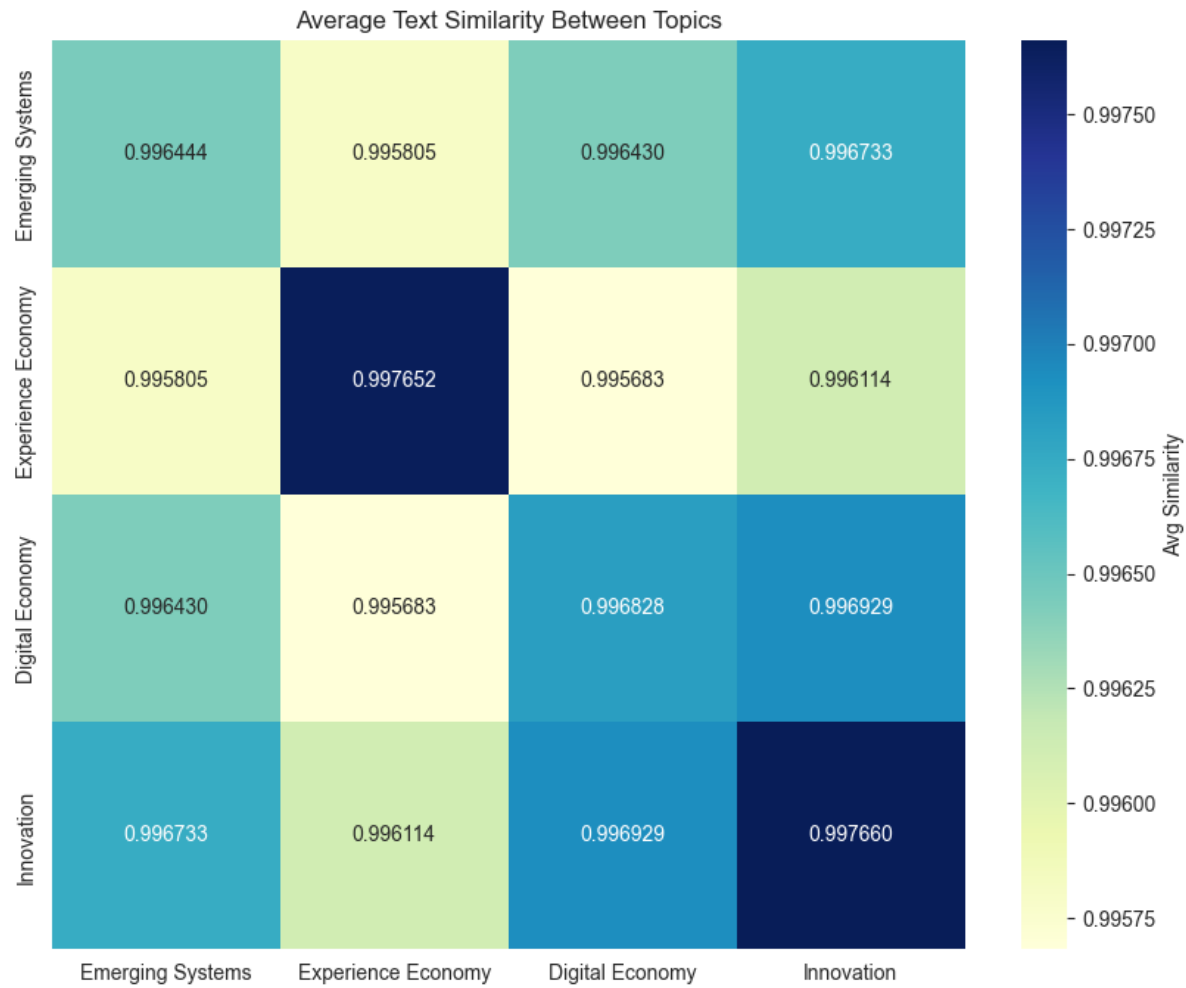


FIGURE 5 SIMILARITY HEAT MAP

	Emerging Systems	Experience Economy	Digital Economy	Innovation
Emerging Systems	0.996444	0.995805	0.99643	0.996733
Experience Economy	0.995805	0.997652	0.995683	0.996114
Digital Economy	0.99643	0.995683	0.996828	0.996929
Innovation	0.996733	0.996114	0.996929	0.99766

$$avg(similarity) = 0.997146$$

IT'S SEEN THAT WITH A SIMILARITY OF 0.997 IT HAS HIGH PROBABILITY OF HAVING EQUAL TOPIC.

ABOUT THIS DOCUMENT

CC BYNCSA. SOME RIGHTS RESERVED.

LICENSE

CREATIVE COMMONS. ATTRIBUTION — NON-COMMERCIAL — SHARE ALIKE 4.0 INTERNATIONAL. FOR FURTHER INFORMATION ABOUT RIGHTS AND RESTRICTIONS, PLEASE VISIT [HTTPS://CREATIVECOMMONS.ORG/LICENSES/BY-NC-SA/4.0/](https://creativecommons.org/licenses/by-nc-sa/4.0/).

THE INFORMATION PROVIDED IN THIS DOCUMENT WAS DEVELOPED PURELY IN PURSUIT OF ACADEMICS FOR UTN FRBA, NATURAL LANGUAGE PROCESSING COURSE, K3551. AS IT WAS TAUGHT IN ITS QUARTERLY FORMAT, BY ENG. HERNAN BORRE.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

AUTHOR

TOMAS AGUSTIN SANCHEZ. FOR MORE INFORMATION ABOUT THE AUTHOR, VISIT [HTTPS://TOMSANCHEZ.COM.AR](https://tomsanchez.com.ar).

REFERENCES

ANDREAS C. MÜLLER, SARAH GUIDO (2016). INTRODUCTION TO MACHINE LEARNING WITH PYTHON: A GUIDE FOR DATA SCIENTISTS. O'REILLY EDITIONS.

HARRY J.W PERCIVAL, BOB GREGORY (2020). ARCHITECTURE PATTERNS WITH PYTHON: ENABLING TEST-DRIVEN DEVELOPMENT, DOMAIN-DRIVEN DESIGN, AND EVENT-DRIVEN MICROSERVICES. O'REILLY EDITIONS.

MATT HANNIBAL (2015). SPACY. OFFICIAL DOCUMENTATION WEBSITE.

SUSAN LI (2018). MACHINE LEARNING FOR TEXT CLASSIFICATION USING SPACY IN PYTHON. TOWARDS DATA SCIENCE MEDIUM PUBLICATION.

SEBASTIÁN RAMÍREZ MONTAÑO (2023). FASTAPI. OFFICIAL DOCUMENTATION WEBSITE.

APENDIX A

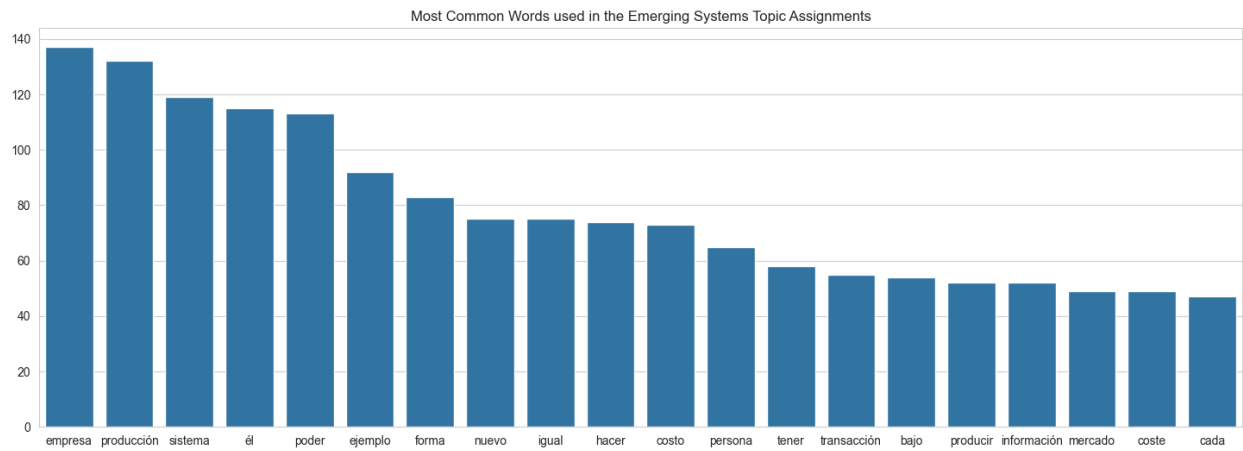


FIGURE 6 MOST COMMON WORDS IN EMERGING SYSTEMS

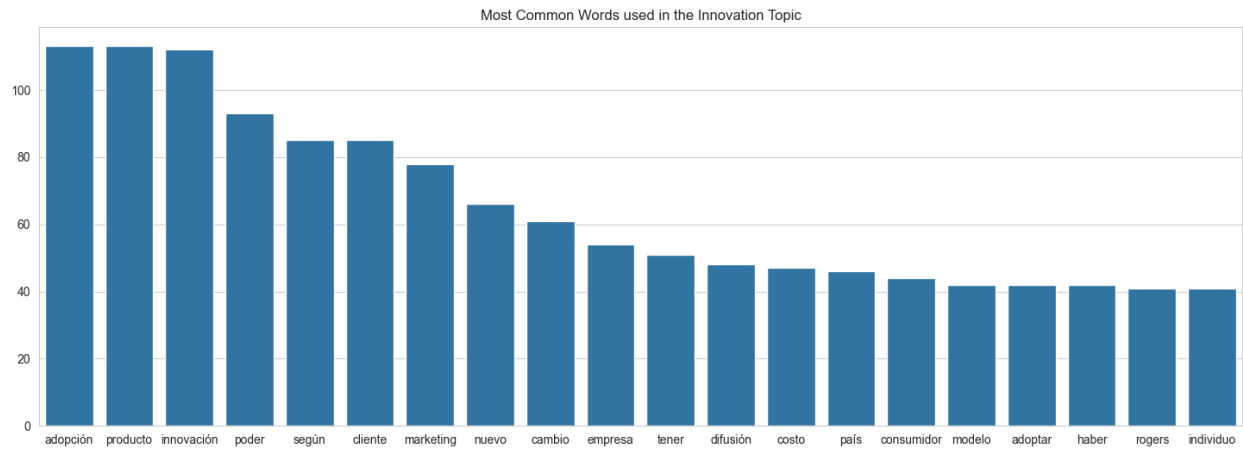


FIGURE 7 MOST COMMON WORDS IN INNOVATION

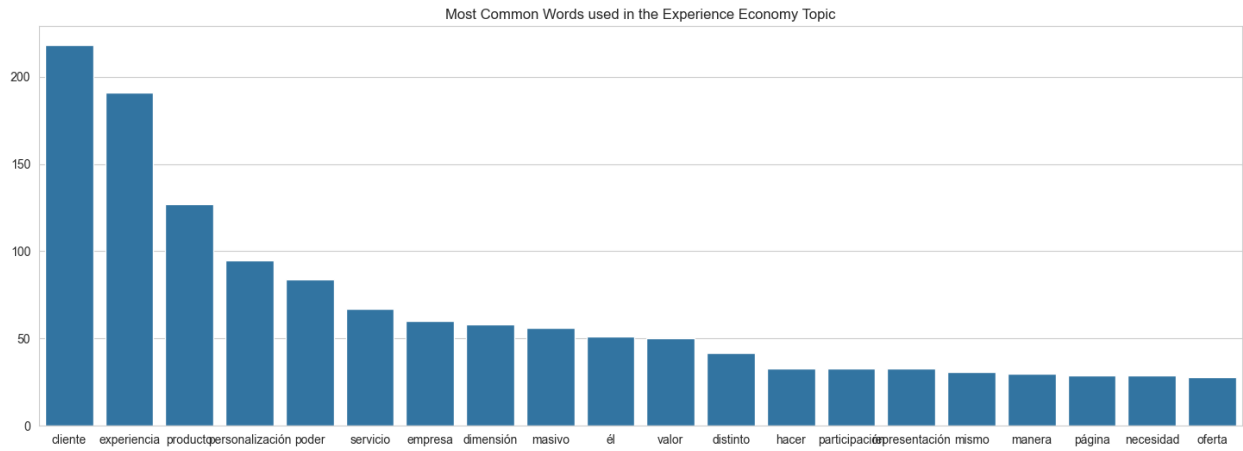


FIGURE 8 MOST COMMON WORDS IN EXPERIENCE ECONOMY