



**Project Title:** Sensing and predictive treatment of frailty and associated co-morbidities using advanced personalized models and advanced interventions

**Contract No:** 690140

**Instrument:** Collaborative Project

**Call identifier:** H2020-PHC-2014-2015

**Topic:** PHC-21-2015: Advancing active and healthy ageing with ICT: Early risk detection and intervention

**Start of project:** 1 January 2016

**Duration:** 36 months

## **Deliverable No: D4.10**

### **LingTester (Prototype) (vers a)**

**Due date of deliverable:** M12 (31<sup>th</sup> December 2016)

**Actual submission date:** 31<sup>th</sup> December 2016

**Version:** 1.3

**Date:** 22<sup>nd</sup> December 2016

**Lead Author:** Sgarbas Kyriakos (UoP)

**Lead partners:** UoP



Horizon 2020  
European Union funding  
for Research & Innovation

<b>Ver.</b>	<b>Date</b>	<b>Status</b>	<b>Author (Beneficiary)</b>	<b>Description</b>
1.0	15/11/16	Draft	N. Fazakis, C. Tsimpouris, K. Sgarbas (UoP)	First Draft
1.1	22/11/16	Draft	N. Fazakis, C. Tsimpouris	Update of introduction, list of tables, minor text corrections, addition of references
1.2	21/12/16	Pre-Final	N. Fazakis, C. Tsimpouris, K. Sgarbas (UoP)	Accompanying files included
1.3	22/12/16	Final	N. Fazakis, C. Tsimpouris, K. Sgarbas (UoP)	Final version

## EXECUTIVE SUMMARY

LingTester is the FrailSafe language analysis tool that aims to process the user's typed text and detect abnormal behaviour. At this point, the prototype is in early alpha stage, but still it is able to perform classification according to levels of frailty. The present deliverable describes the development of the prototype, the algorithms used, the training process and some preliminary test results.

This deliverable is part of WP4. The main objective of this Work Package is to handle the collection, management and analysis of frailty older people data streamed through their social, behavioural, cognitive and physical activities. Both offline and online methods will be developed. Moreover, the above methods will be applied in order to manage and analyze new data and also generate the FrailSafe patient models.

LingTester will be able to detect signs of mental frailty and personality trait shifts by linguistic processing of a person's written (typed) messages. The linguistic analysis is performed in several layers (ranging from word spelling to Part of Speech -POS- analysis) utilizing the state of the art models in order to determine the mental states of the patients the input texts exhibit. The linguistic corpus obtained from D4.7 is used both for the initial training and the final passive mode (off-line) testing of the prototype.

Along with the results of the patient data analysis, two aiding software programs are delivered in order for the FrailSafe user to manage the patient database and use prediction model to obtain predictions for specific potential patients.

**DOCUMENT INFORMATION**

<b>Contract Number:</b>	H2020-PHC–690140	<b>Acronym:</b>	FRAILSAFE
<b>Full title</b>	Sensing and predictive treatment of frailty and associated co-morbidities using advanced personalized models and advanced interventions		
<b>Project URL</b>	<a href="http://frailsafe-project.eu/">http://frailsafe-project.eu/</a>		
<b>EU Project officer</b>	Mr. Jan Komarek		

<b>Deliverable number:</b>	4.10	<b>Title:</b>	LingTester (Prototype) (vers a)
<b>Work package number:</b>	4	<b>Title:</b>	Data Management and Analytics

<b>Date of delivery</b>	<b>Contractual</b>	<b>31/12/2016 (M12)</b>	<b>Actual</b>	31/12/2016
<b>Status</b>	Draft <input checked="" type="checkbox"/>		Final <input type="checkbox"/>	
<b>Nature</b>	Report <input type="checkbox"/>	Demonstrator <input type="checkbox"/>	Other <input checked="" type="checkbox"/>	
<b>Dissemination Level</b>	Public <input checked="" type="checkbox"/>	Consortium <input type="checkbox"/>		
<b>Abstract</b> (for dissemination)	This is a public deliverable that summarizes the progress of the construction of LingTester offline prototype. The architecture of the system is described in detail, and what steps were needed to evaluate its output. Also, the structure of the internal offline database is described and how it is managed. The deliverable includes a demo and source files.			
<b>Keywords</b>	offline data management, frailty prediction, classification			

<b>Contributing authors</b> (beneficiaries)	Fazakis Nikos (UoP) Tsimpouris Charalampos(UoP) Sgarbas Kyriakos (UoP)			
<b>Responsible author(s)</b>	Sgarbas Kyriakos	<b>Email</b>	sgarbas@upatras.gr	
	<b>Beneficiary</b>	UoP	<b>Phone</b>	+30 2610 996 470

**Table of Contents**

**CHANGE HISTORY .....2**

**EXECUTIVE SUMMARY .....3**

**DOCUMENT INFORMATION.....4**

**List of Figures .....7**

**List of Tables .....8**

**List of abbreviations and acronyms (*in alphabetical order*) .....9**

**1. Introduction..... 10**

**2. LingTester architecture ..... 11**

    2.1 Initial architecture ..... 11

    2.2 Current architecture ..... 12

**3. Data collection ..... 13**

    3.1 Data analysis ..... 14

    3.2 Data verification ..... 17

**4. Local database ..... 18**

    4.1 Introduction ..... 18

    4.2 Database description ..... 18

    4.3 Part-Of-Speech extraction..... 21

    4.4 English translation for sentiment analysis ..... 22

    4.5 Database self-validation ..... 23

    4.6 Database auto-update..... 23

**5. Experiments and results..... 23**

    5.1 Preprocessing ..... 23

    5.2 Precision & Recall ..... 26

**6. Conclusions ..... 32**

**7. System development ..... 33**

    7.1 Software used for training and prediction..... 33

    7.2 Installation notes for training and prediction..... 33

        7.2.1 Python core and Python modules..... 34

        7.2.2 Java Virtual Machine (JVM)..... 34

        7.2.3 POS Tagger..... 34

        7.2.4 OpenOffice dictionaries ..... 34

        7.2.5 Weka ..... 34

**8. Frailsafe user software ..... 35**

    8.1 Steps to import new data ..... 35

<b>FRAILSAFE – H2020-PHC–690140</b>	<b><i>D4.10</i></b>
8.2 Steps to export results .....	36
<b>9. Future improvements for the Prototype.....</b>	<b>39</b>
<b>10. Ethics and Safety .....</b>	<b>40</b>
<b>11. References .....</b>	<b>41</b>
<b>12. Source files .....</b>	<b>43</b>
<b>13. Annexes .....</b>	<b>44</b>
13.1 Database management & feature extraction.....	44
13.2 Prediction tool .....	55

## List of Figures

<a href="#">Figure 1. Initial architecture for LingTester</a>	11
<a href="#">Figure 2. Current architecture for offline LingTester</a>	13
<a href="#">Figure 3. Patients per language</a>	14
<a href="#">Figure 4. Patients per frailty status</a>	15
<a href="#">Figure 5. Distribution of patients per sex</a>	16
<a href="#">Figure 6. Patient distribution per transcript feature</a>	24
<a href="#">Figure 7. Training and prediction methodology</a>	24
<a href="#">Figure 8. Decision tree based on C4.5 algorithm</a>	27
<a href="#">Figure 9. Model statistics</a>	30
<a href="#">Figure 10. Screenshot of real time operation</a>	36
<a href="#">Figure 11. Input file structure</a>	37
<a href="#">Figure 12. Screenshot of prediction programme while executing</a>	38

## List of Tables

<a href="#">Table 1. POS tagger example</a>	22
<a href="#">Table 2. Analysis of all extracted features</a>	25
<a href="#">Table 3. Summary of classification algorithms</a>	29
<a href="#">Table 4. Parameterization of the Decision Tree model</a>	29
<a href="#">Table 5. Parameterization values of the Decision Tree model</a>	30

## List of annexes

Annex 1	<a href="#">Database management &amp; feature extraction</a>	44
Annex 2	<a href="#">Prediction model, source code</a>	55



**List of abbreviations and acronyms** *(in alphabetical order)*

API	Application Programming Interface
ARFF	Attribute-Relation File Format
JVM	Java Virtual Machine
KNN	K-Nearest Neighbor
NLP	Natural Language Programming
LOOCV	Leave-one-out cross-validation
PoS	Part Of Speech
SVM	Support Vector Machine
UoP	University of Patras
Weka	Waikato Environment for Knowledge Analysis

# 1. Introduction

A proper evaluation of the nature of a patient's language impairment requires consideration of patterns of breakdown in the context of an account of language comprehension which specifies the various processes and representations involved: the representation of linguistic knowledge, its automatic and controlled access, and the mental processes which combine different types of linguistic knowledge. When we understand a written sentence, we automatically access the meanings of individual words together with their syntactic specifications and combine them on the basis of this lexically specified information. There are various types of text analysis processes; some are syntactic, such as the integration of a definite article with the following noun or adjective to form a noun phrase, and the integration of a verb and its argument into a verb phrase. Others are morphological, such as the combination of stems and affixes to form morphologically complex words, and yet others involve combinatorial processes which modify the meanings of words when they are integrated with other words. For example, an aspect of the meaning of grass is that it is green; but in the phrase dry grass, its meaning changes slightly to highlight a different aspect of grass— its brownness.

In order to implement the first version of LingTester an architectural analysis was prepared. The current architecture of LingTester led the need to categorize the research in four main areas.

The first area of research relates to the analysis and the issues of the collected subjects' data. The key aspects and features of the dataset were investigated and solutions were given to the unexpected arisen problems. The structural organization of the local database and the development of the processes for the management of the local database was another area of research. A simple but very informative and mobile structure for offline data storage along with its necessary manipulation methods was designed and implemented. In the area of the highly critical classification task, the domains of feature extraction, feature selection text classification and model optimization were deeply studied and exhaustively experimented in order to obtain the first acceptable prediction results. The implementation of the aiding software for the FrailSafe user was another key area of research, with the deployment of technologies like Python and Java a cross-platform approach was achieved and the first semi-integrated user software package has been developed successfully.

Following the analysis of the basic areas of research this deliverable includes, chapters relating with the System Development subjects of the software used and the installation notes, the basic usage instructions of the FrailSafe software package and finally the future work that seems promising and is already being planned, for the completeness of the document.

## 2. LingTester architecture

### 2.1 Initial architecture

The initial design of the Natural Language Analysis component (a.k.a. LingTester), is shown in the following diagram:

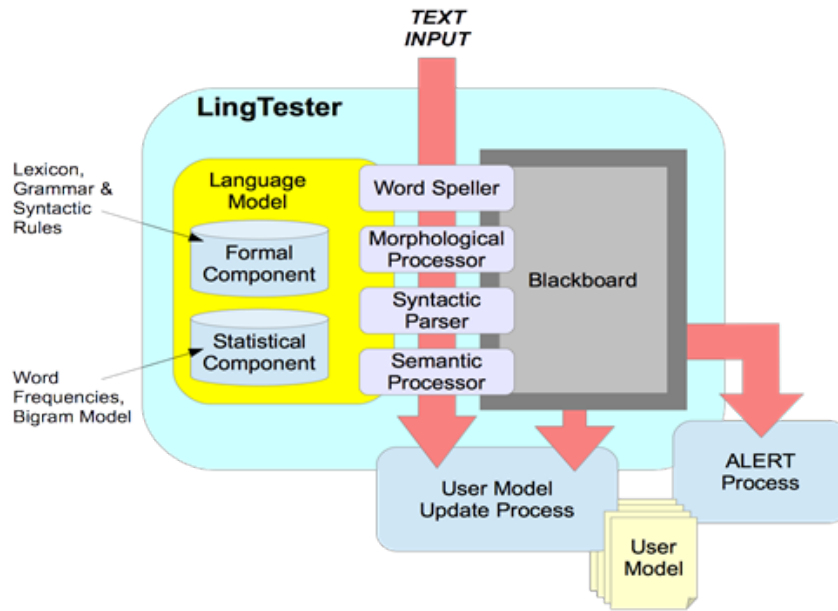


Figure 1. Initial architecture for LingTester

LingTester was initially designed to include four computational linguistic modules: a Word Speller, a Morphological Processor, a Syntactic Parser, and a Semantic Processor. A Language Model would be able to feed them with linguistic information. The Language Model would be composed by a Formal Component representing the lexicon, grammar and syntactic rules of the language, and a Statistical Component containing results of word and bigram frequencies. The whole structure was meant to be modular (in order to facilitate its use in several languages - two in the project) and would be developed over a blackboard scheme representation model that would be able to collect the output of each component and enable them to interact with each other, when necessary. The system was based on the main assumption that all texts would be written by the patients themselves and the classification module would provide rule-based results.

Thus, after a stage of adaptation/training/fine-tuning, LingTester was *expected* to detect frailty symptoms related to the use of language, and derive the patient's condition, according to the following table:

STAGE	CONDITION	SYMPTOMS (Indicative)
1	Normal	-
2	Cognitive Decline	Misspellings, character transpositions/ eliminations.
3	Cognitive Impairment	Misuse of functional words, morphological errors.
4	Dementia	Serious syntactic and semantic errors.

It should be stressed once more that the aforementioned design presupposes that the users will be able to type their own text messages, since key features like spelling, morphology and punctuation carry a significant portion of the information we expected to use in order to detect possible discrepancies in written text.

With this in mind, we have included some questions in the questionnaires (detailed in deliverable D2.1 Clinical Study Methodology) that serve the purpose of producing data for the LingTester (for training and test). However, during the process of questionnaire collection we realized that very few subjects were able to type. The vast majority were not, and they dictated their answers to the person responsible for carrying out the interview who then typed the answers as best as they could.

This changed the initial architecture significantly. We had to use methods of mental frailty detection that would be robust under interpretation via a third person. Specifically, it became evident to treat the process of text writing as a Hidden-Markov Process, where the person's mental state evolves over time, but we do not have direct data on this evolution, but only indirect evidence (the written text) that can reveal the (hidden) mental state. For this reason, statistical and information-based text analysis methods were preferred instead of the rule-based of the initial design. This was already documented in D2.1. The new (current) architecture is described next.

## 2.2 Current architecture

After thorough study of current practices and development, we re-evaluated the initial architecture according to the following figure. The structure has been kept mostly intact, however with some specific changes. Firstly, the syntactic parser has been replaced by a Part Of Speech tagger (see section 4.3 for more details), as the later can produce more accurate results for the Greek language. Furthermore, semantic processing module was replaced by a sentiment analysis module to try and provide a sentiment analysis of the patient through written text. Written texts are translated to English one, and then we can export polarity features to the

classification process. This decision was based at the fact, that there is no state of the art available tools for sentiment analysis for this language. It has been shown that working with standard technology and existing sentiment analysis approaches is a viable approach to sentiment analysis within a multilingual framework (Denecke, 2008).

As shown in the following figure (Figure 2), written text is submitted to LingTester tool through a predetermined process and is stored within a secure database for further analysis. In order to create the training model, all patient rows are fetched from the offline database and features are extracted for the next step. Each feature utilises different resources and is based on different third-party or not tools. These tools are described thoroughly in section 4. This step is followed by the training module which extracts a model in a binary format for testing and evaluation. This methodology has been repeated multiple times so as to maximise accuracy while also optimising all parameters of the system. The final model is packaged in a way to be more programmer-friendly (see section 8 for more details).

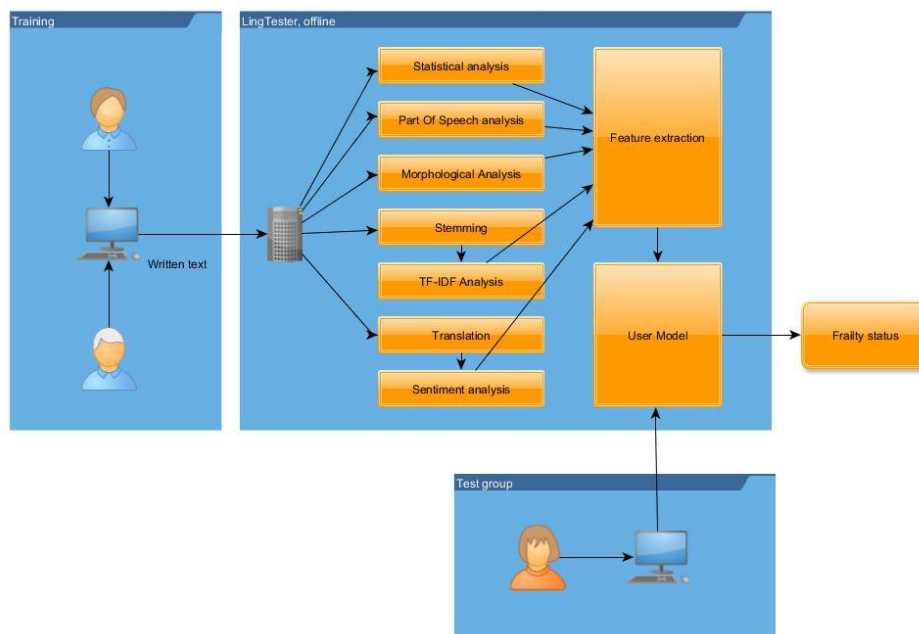


Figure 2. Current architecture for offline LingTester

Finally, it should be stressed that system should provide patient's condition in the following clusters: non-frail, pre-frail, frail.

### 3. Data collection

Although data collection at first might not seem related to the development of the prototype, it actually plays a very significant role. Even inspection of the data collected can indicate which analysis tools and methods are applicable to the task. We have already stated that the nature of the collected data forced us to reconsider the aforementioned initial design. Moreover, in the

current architecture text data are needed to train and test several components of LingTester. For this first version of the prototype, the texts we used were the ones that were available until 31st October 2016. This first batch of text data is described next.

### 3.1 Data analysis

The current deliverable uses data collected until 31/10/2016. Until then we had available the following patient data:

- Data from 51 patients, UoP, Greece. One submission per patient.
- Data from 66 patients, MATERIA, Cyprus. One submission per patient.

The following patient data per recruitment centre are expected, according to the following timeline as set in detail at *D2.1 Clinical study methodology*:

- 80 patients from Start-up Group A,
- 20 patients from main Group B,
- 25 following afterwards will belong to the Evaluation Group C and
- the last 25 to the Control Group (D), totalling 150 patients.

The total data volume from all recruitment centers will include 450 patients, with multiple submissions during the duration of the project, per patient. An initial statistical analysis based to the current set of patient data returned the following results per classified feature:

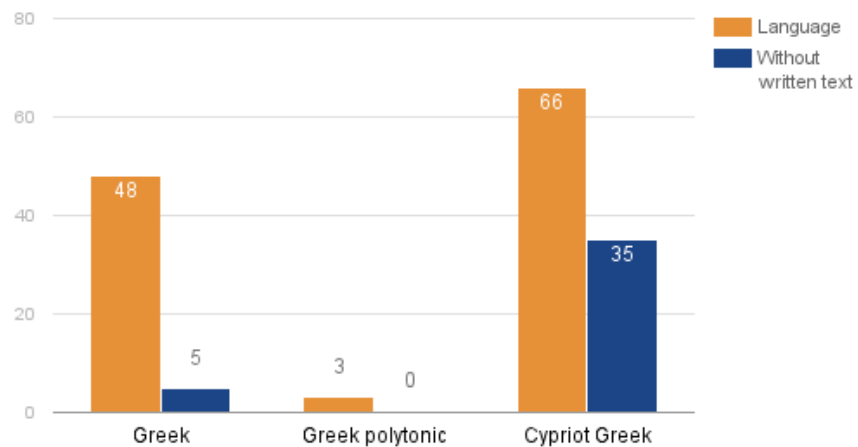


Figure 3. Patients per language

- Per language:

- Data from 48 patients were provided in Greek, but 5 of them refused or were unable to provide any written text, neither for the description of image, nor describing a personal event
- Data from 3 patients were written in Greek polytonic. While it was not in our initial scope to differentiate this information, it was recorded in the offline database for future study
- Data from 66 patients were provided in Cypriot Greek, but 35 of them refused or were unable to provide any written text, neither for the description of image, nor describing a personal event

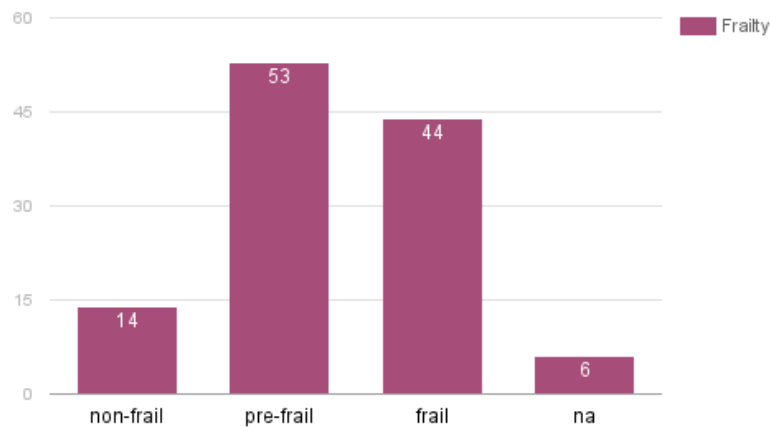


Figure 4. Patients per frailty status

- Per frailty status:
  - 14 patients were classified as non-frail
  - 53 patients were classified as pre-frail
  - and 44 patients were classified as frail
  - while for 6 patients, there was no information available for their frailty status. At this point, we should stress that all these patient data (6) were excluded from the training procedure

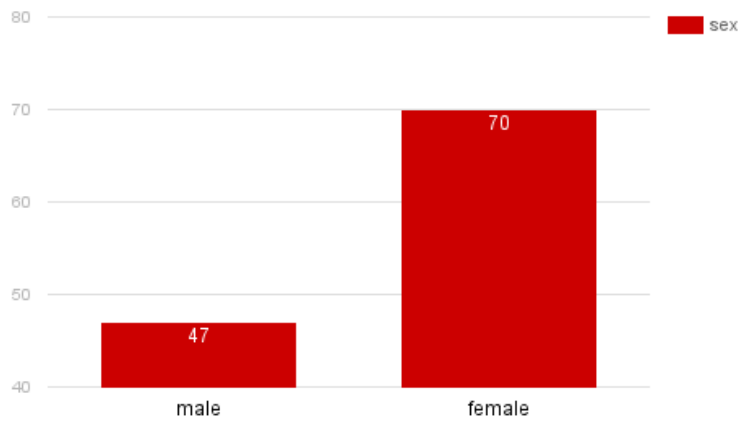


Figure 5. Distribution of patients per sex

- Per sex:
  - 47 patients were male
  - and 70 patients were female

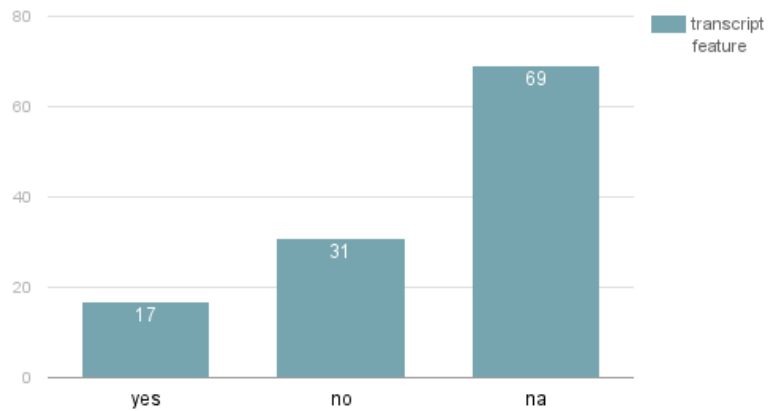


Figure 6. Patient distribution per transcript feature

- Per transcript mode:
  - 17 transcripts (column “yes”) were dictated from the patient to the doctor, and the latter typed the text



- 31 transcripts (column “no”) were handwritten by the patient, and they were afterwards digitized by the doctor
- 69 was classified as na (not available) meaning that patient refused to participate in this action (written or oral) which are 40 in total, see figure 3 column “without written text”. Furthermore, for the remaining 29 cases we were unaware of which case it was between yes/no and this is the reason we set this attribute to na, till we have a detailed report of which is the case by the providing team. We project, that this number will be lower in the revised report.

Finally, in total only two patients provided written text outside the aspect if the project. As this number is extremely small, we decided not to include this data in the following research methodology.

## 3.2 Data verification

During data collection, some discrepancies were noted concerning the slightly different patient numbering among different research groups. In the local database we took measures to verify the correct index for each patient since no version control was used during the collection of each dataset.

## 4. Local database

### 4.1 Introduction

An offline database has been created based on initial raw data, as given by the partners. While we are still under heavy development till we finalise the selection of features, we had to keep an easy-to-use, cross-platform or operating system, and loose structure database. This aforementioned decision provided numerous advantages, as discussed below:

- Each patient data is saved within a single text file, with UTF8 encoding, based on the underlying file structure system, whatever this is (NTFS/Fat32 for Windows, Ext4 for Linux, etc). This also helps to avoid any database management tools like ODBC drivers and so forth.
- File naming is based on patient id. So, it is really easy to retrieve data for a single patient and edit the file whenever necessary through user's favourite editor
- All tags (attributes) per patient are dash prefixed, so we can add or remove attributes however it suits best
- Data is retrieved and saved through generic Python functions, as described in the following section. Also, due to the use of the filesystem, we can also construct different or new methodology in a different programming language without potential connection problems with a database system.
- Create backups of all data, using a compressed data type, like zip
- Can support versioning. Using any known version control system such as GIT, we are able to keep track of the database changes at all times looking backwards

On the downside, this database structure does not provide any security firewall by itself. Security is based on the access provided by the file system, and for this reason all files are stored locally.

While it can be argued if this this structure can sustain a production-ready solution with thousands of rows, we keep in mind at all times. Upon methodology finalisation, database will also get its final form and this will be discussed again in a later report.

### 4.2 Database description

The structure of the data within the database files is described below. An example of a patient's data with id 1001 is the following:

**-patient** 1001  
**-tag** prefrail  
**-transcript** no  
**-language** greek  
**-sex** female

**-desc\_event**

Υπάρχουν πολλά σημαντικά θετικά γεγονότα στη ζωή μου. Ένα από αυτά είναι η επιτυχία μου στη Φιλοσοφική Σχολή του Πανεπιστημίου Αθηνών, καθώς την ίδια περίοδο, με παράλληλες εξετάσεις και στη Νομική Αθηνών. Πήρα μεγάλη χαρά γιατί ήταν καρπός πολύ εντατικής μελέτης- ένα ολόκληρο καλοκαίρι μόνον- με δεκαεπτά ημέρες μόνον φροντιστήριο και επί πλέον είχα μεγάλη επιθυμία να σπουδάσω, να μορφωθώ. Διάλεξα τη Φιλοσοφική και δεν έχω ούτε μια στιγμή μετανοιώσει...

**-desc\_event\_ENG**

There are many important positive events in my life. One of them is my success at the Philosophical Faculty of the University of Athens, as the same period, with parallel tests and Athens Law. I took great pleasure because it was the result of very intensive meletis- an entire summer monon- with seventeen days only tutorial and moreover I had a great desire to study, get an education. I chose Philosophy and I have not a moment regret ...

**-desc\_event\_POS**

Υπάρχουν verb/--/active/plural/present  
 πολλά adjective/accusative/neuter/plural/--  
 σημαντικά adjective/accusative/neuter/plural/--  
 θετικά adjective/accusative/neuter/plural/--  
 .  
 .  
 .  
 έχω verb/--/active/singular/present  
 ούτε conjunction/--/--/--/--  
 μια numeral/--/--/--/--  
 στιγμή noun/accusative/feminine/singular/--  
 μετανοιώσει... noun/genitive/masculine/singular/--

**-desc\_image**

Βρισκόμαστε μπροστά σε μια ‘ σουρεαλιστική’ εικόνα, σε μια σκηνή που διαδραματίζεται σε μια κουζίνα. Μια ‘καλή’ νοικοκυρά ασχολείται με το πλύσιμο ή το σκούπισμα των πιάτων ενώ μπροστά της η λεκάνη του νεροχύτη πλυμμηρίζει και τα νερά χύνονται στο πάτωμα. Θαυμάζει κανείς τη μακαριότητά της, την αταραξία της μπροστά στο φαινόμενο. Άραγε τί την απασχολεί που δεν μπορεί να αντιληφθεί ότι τα παιδιά της, στον ίδιο χώρο, πίσω από την πλάτη της... ‘ κλέβουν’ το γλυκό από το επάνω ντουλάπι και το χειρότερο, ο γιός της που έχει ανέβει πάνω σε ένα ψηλό σκαμπό κοντεύει να πέσει, καθώς αυτό γέρνει στο πλάϊ ,έτοιμο να καταρρεύσει. Ω, τι κόσμος μαμά!!

**-desc\_image\_ENG**

We face a surreal picture, in a scene that takes place in a kitchen. A good housewife engaged in washing or wiping the dishes while in front of the sink basin plymmirizei and the water poured on the floor. One admires the blessedness, the equanimity of the front of the phenomenon. I wonder what the concern can not perceive that children of the same place, behind the ... back. Steal sweet from the upper cabinet and the worst, the son who has climbed on a tall stool is about to fall, as it leans on the side, ready to collapse. Oh, mama world !!

```

-desc_image_POS
Επισκόμαστε verb/--/active/plural/present
μπροστά adverb/--/--/--/--
σε preposition/--/--/--/--
μια numeral/--/--/--/--
.
.
.
τι pronoun/inflectionless/--/--/--
κόσμος noun/nominative/masculine/singular/--

-prev_text

-prev_text_ENG

-prev_text_POS
    
```

While, there is no styling within the plain text files, and newlines do not affect the parsing of the file structure, from the above structure we can easily retrieve all available attributes (tags) as they are all prefixed by a dash, which are:

- **-patient:** The patient ID. This attribute exists in all files. While the same number exists also in the filename, we put it here for consistency reasons and backwards compatibility for future updates.
- **-transcript:** This is identified by the following options (also described in detail in section 3.1)
  - yes: Text was written by the doctor, while the patient was talking
  - no: Text was written in hand by the patient, and it was digitized through the doctor
  - na: Not available, for instance for patients that refused to participate in this action
- **-language:** This is identified by the following options
  - greek: Text is in Greek
  - greek-polytonic: Text is in Greek Polytonic.
  - greek-cypriot: Text is in Greek cypriot.
  - french: Text is in French.
- **-tag:** This is identified by the following options
  - nonfrail: Patient is identified as non-frail
  - prefrail: Patient is identified as pre-frail
  - frail: Patient is identified as frail
  - na: Data is missing/not available
- **-sex:** This is identified by the following self-explanatory options
  - male
  - female
- **-desc\_event:** Multiline text, the description of an event
- **-desc\_image:** Multiline text, the description of an image

- **-prev\_text**: Multiline text, previous text of the same patient, which is not necessarily a description of an event or an image. It can be of any context and is provided by the subject, for instance an old email, to compare extracted features between different time periods.
- **-desc\_event\_POS, -desc\_image\_POS, -prev\_text\_POS**: Part of speech information for each multiline tag, as set before
- **-desc\_event\_ENG, -desc\_image\_ENG, -prev\_text\_ENG**: English translation, based on **-desc\_event, -desc\_image, -prev\_text** data.

### 4.3 Part-Of-Speech extraction

The Part of Speech tagger attempts to automatically determine the part of speech (e.g., noun, adjective, verb, etc.) of each word occurrence in Greek texts. It can also tag each word occurrence with additional information, such as the gender, number, and case of each noun, the voice, tense, and number of each verb (Koleli, 2011). The current version of AUEB's Greek POS tagger that was used is version 2 alpha and is released under the GNU General Public License.

The POS tagger can recognise the following classes of words, along with other useful information per case:

1. adjective
2. adverb
3. article
4. conjunction
5. noun
6. numeral
7. other
8. particle
9. preposition
10. pronoun
11. punctuation
12. verb

Following is an example evaluating this POS tagger in action. The sentence “Υπάρχουν πολλά σημαντικά θετικά γεγονότα στη ζωή μου.” (which translates to “*There are many important positive facts in my life.*”) produces the following information.

<b>Υπάρχουν</b> <i>There are</i>	verb	--	active	plural	present
-------------------------------------	------	----	--------	--------	---------

πολλά <i>many</i>	adjective	accusative	neuter	plural	-
σημαντικά <i>important</i>	adjective	accusative	neuter	plural	-
θετικά <i>positive</i>	adjective	accusative	neuter	plural	-
γεγονότα <i>facts</i>	noun	accusative	neuter	plural	-
στη <i>in</i>	article	prepositional	accusative	feminine	singular
ζωή <i>life</i>	noun	accusative	feminine	singular	-
μου <i>my</i>	noun	genitive	masculine	singular	-
.	punctuation	-	-	-	-

Table 1. POS tagger example

In order to optimize the system, POS information is extracted once per case, and the database is automatically updated for future use. See also section 4.6 for more details.

## 4.4 English translation for sentiment analysis

To achieve better results in sentiment analysis, a significant decision was made to avoid direct sentiment analysis in Greek language or French one, but translate texts in English and then evaluate the later one. This process, also helps make the system even more language independent by utilising a unified translation system, and then shifting the sentiment analysis problem to a different level.

However, in order for this methodology to work, a third party translation service had to be used. After further investigation, we narrowed down to MyMemory service (<https://mymemory.translated.net/>), a **free to use** service. This translation service, uses both human and machine learning techniques for best results. MyMemory gives quick access to a large number of translations originating from professional translators, LSPs, customers and multilingual web content. It uses a powerful matching algorithm to provide the best translations available for the source text. Last but not least, we should mention that MyMemory currently contains professionally translated segments. System is constructed in a way to be modular in

mind, and this is also the case for the translation submodule. In case there is any discontinuance of this service, we can easily switch to a different one, like the well known paid Google Translation API service. However, MyMerory was used for its simple API, free of charge pricing while providing translations of high quality results.

## 4.5 Database self-validation

As database management is one of the many steps to extract and create the frailty status prediction model, we had to be sure that all data stored should fulfill all the requirements for the next steps of analysis, which is data/text mining and will be performed by the free software WEKA (Eibe, Mark, Ian 20016). Having said that, it was of paramount importance for the created ARFF files (the WEKA-specific file format) to be always valid, avoid missing attributes or identify typos. So, a function was created for this purpose that reads all patient data and tries to identify discrepancies between saved data and expected classes. Finally, it also exports some basic statistical information, as were explained in detail in section 3.1

## 4.6 Database auto-update

In order to have a valid state of the database at all times, a library has been constructed in a way to always keep its internal structure useable and filled with all needed data. For this reason, specific functions have been created to check and update all patient data according to user standards. Function `update_pos_info_everywhere()` loads all patient data available from the database, and updates any missing POS data for new or updated rows. Also, function `update_corpus_with_translations()`, whenever called, will update the whole database with any missing english translations. The function has been constructed in a way to respect initial language as set in each patient data, which means that translation respects patient's language. For example Greek to English for Greek patients and French to English for patients from France and so forth.

# 5. Experiments and results

## 5.1 Preprocessing

### *Feature extraction*

The classification task of the mental state of a subject requires the deployment of machine learning and pattern recognition techniques. The basic requirement for these techniques is the

processing of the organized patient data with feature extraction methods before the training and prediction procedures as can be viewed in [fig. 7](#). Feature extraction starts from an initial set of measured data and builds derived values (features) intended to be informative and nonredundant, facilitating the subsequent learning and generalization steps, and in some cases leading to better human interpretations. Feature extraction involves reducing the amount of resources required to describe a large set of data. When performing analysis of complex data one of the major problems stems from the number of variables involved. Analysis with a large number of variables generally requires a large amount of memory and computation power, also it may cause a classification algorithm to overfit to training samples and generalize poorly to new samples. Feature extraction is a general term for methods of constructing combinations of the variables to get around these problems while still describing the data with sufficient accuracy.

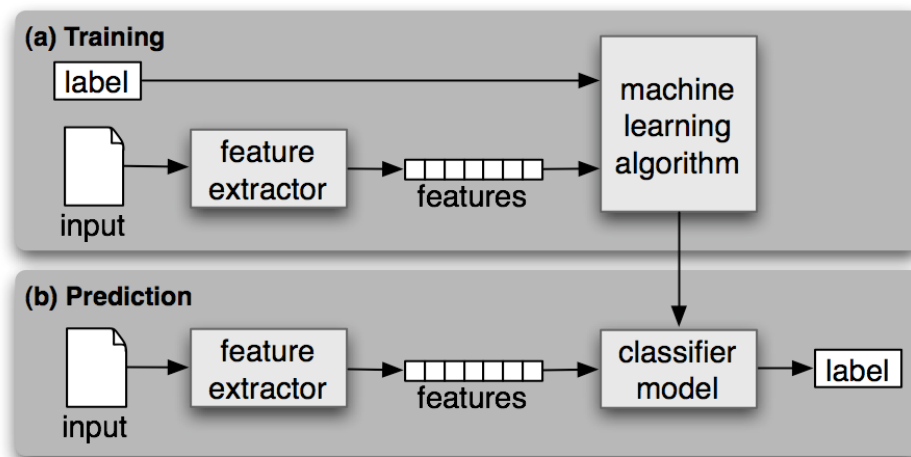


Figure 7. Training and prediction process

The implemented feature extraction algorithm for the LingTester uses several extraction methods (see [table 3](#)). The first ones involve the standardization of the basic attributes of the collected data. For the features transcript, language, class and sex, which as their name suggests describe basic information from collected data, simple rules and correction algorithms have been applied in order for the extracted data to be distinctly formalized.

Another categorization of feature extraction methods implemented uses statistical measures for the written text of the subjects. Those measures include the text length, the number of sentences, the number of words per sentence and the text entropy. The text entropy, a measure of unpredictability of information content based on characters.

Proceeding to more NLP specific techniques the term frequency–inverse document frequency (tf-idf) is used. Tf-idf is a numerical statistic that is intended to reflect how important a word is to a document in a corpus. It is used as a weighting factor in text mining. The tf-idf value increases proportionally to the number of times a word appears in the document, but is offset by the



frequency of the word in the corpus, which helps to adjust for the fact that some words appear more frequently in general (Salton et al, 1983). To gain as much information as possible from this methodology, we utilised tf-idf twice. The first time is based in stemmed words, in order to avoid all suffixes. The second one, is based on POS data. This way, we could identify possible unigrams, or bigrams that are more frequent than other (for instance verb+adjective).

Written text can be broadly categorized into two types: facts and opinions. Opinions carry people's sentiments, appraisals and feelings toward the world. The module(open source) that is used for sentiment analysis (`sentiment` within `pattern.en`) bundles a lexicon of adjectives (e.g., good, bad, amazing, irritating, ...) that occur frequently in product reviews, annotated with scores for sentiment polarity (positive ↔ negative) and subjectivity (objective ↔ subjective). Using the `sentiment()` function we gain polarity and subjectivity for the given sentence, based on the adjectives it contains, where polarity is a value between -1.0 and +1.0 and subjectivity between 0.0 and 1.0.

A last preprocess step was to try and identify misspellings. In order to base our work on open source or community based tools, we used the python `pyenchant` library combined with the OpenOffice speller dictionary. This speller, is widely used by thousands of users through OpenOffice applications like OpenOffice Writer, for multiple operating systems like Linux or Microsoft Windows and is easily accessible through a Python API. For our case, we extracted the number of misspelling words against all words per case. The following table summarizes all exported features.

Feature Names	Extraction Method
Transcript, Language, Class, Sex	Rules & Correction filters
Text_length, Number_of_sentences, Number_of_words, Number_of_words_per_sentence, Text_entropy	Statistical Measures
Desc_image_ENG_sentiment, Desc_event_sentiment, Prev_text_ENG_sentiment	Sentiment Analysis
Tf-XX	Term frequency – Inverse document frequency
Tf-pos-XX	Part of Speech analysis, using tf-idf methodology

Table 2. Analysis of all extracted features

## 5.2 Precision & Recall

### *Feature selection*

The next step before proceeding to classification task is the feature selection task. Feature selection, also known as variable selection, attribute selection or variable subset selection, is the process of selecting a subset of relevant features (variables, predictors) for use in model construction. Feature selection techniques are used for three reasons:

- simplification of models to make them easier to interpret by researchers/users,(Gareth, 2013)
- shorter training times,
- enhanced generalization by reducing overfitting(formally, reduction of variance)

The central premise when using a feature selection technique is that the data contains many features that are either redundant or irrelevant, and can thus be removed without incurring much loss of information(Bermingham, 2015). Redundant or irrelevant features are two distinct notions, since one relevant feature may be redundant in the presence of another relevant feature with which it is strongly correlated.

A number of techniques have been proposed in the literature using algorithms and even classifiers for automating the process of feature selection. The most common algorithms are the exhaustive, best first (Pearl, 1984), simulated annealing (Khachaturyan, 1979) and the genetic algorithm (Mitchell, 1996). In practice, the task of feature selection is a highly empirical process where algorithms and human intelligence are combined in order to find the optimal subset of features, thus constructing the final feature set that will be used in the classification task.

As a first approach to feature selection, a simple process has been followed. The first steps of the process involve an iteration of classifications where each individual feature was examined for its contribution to the accuracy of the temporary model, using the cross validation method(Geisser, 1993). After a sufficient number of iterations, the resulting decision tree was visualized and examined by hand in order to further optimize the resulting model.

### **Feature Selection Algorithm**

---

**Input:**

Load the complete set of features (C)

Count the number of all features (N)

Classify with C and store the accuracy (A)

Initialize pointer as zero (P)

**Loop for N**

Remove C[P]

Classify with C (Ac)

If Ac < A

Restore C[P]

**Validate** features by tree visualization

**Output:**

Subset of features (S)

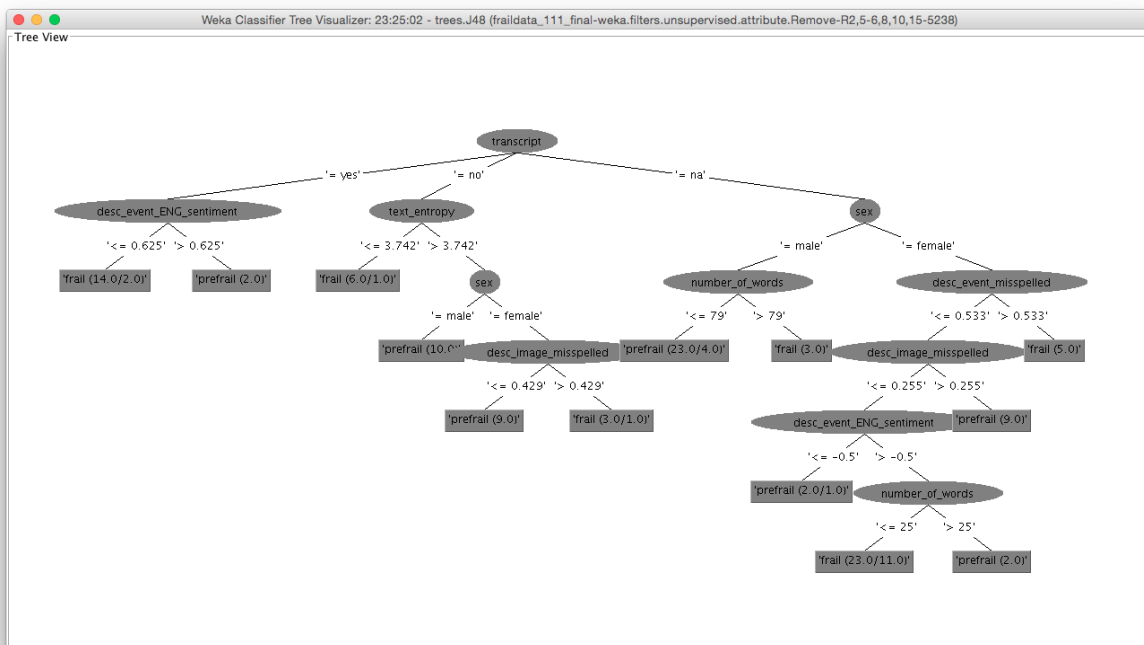


Figure 8. Decision tree based on C4.5 algorithm

The final step of tree visualization was intentionally added in order to exploit a fundamental property of the way decision trees build their structure. In more detail, these models use the information gain (Mitchell, 1997) to rank and place the best contributing features on the top of the tree. Thus, it was possible to validate the importance of its selected feature, understand its contribution and remove the remaining non important features.

**Classification process**

The automatic classification of documents into predefined categories is an important field of active research, the documents can be classified by three classes of methods:

- Unsupervised(Duda, 2001) methods, where no human intervention is required for labeling the collected data and the algorithms deployed are responsible for grouping the data to distinct categories.
- Supervised methods, usually the human expertise is used for labeling each individual instance of the dataset.
- Semi supervised methods, in this class of methods as little as possible human expertise is required to label a small initial amount of data and the algorithms exploit the existence of unlabeled data in order to enrich the training dataset.

The last few years, the task of automatic text classification has been extensively studied and rapid progress seems in this area, the machine learning approaches include the use of classifiers like Bayesian classifier (Russel, 2003), Decision Tree, K-nearest Neighbors (KNN), Support Vector Machines (SVMs) (Cortes, 1995) and Neural Networks (McCulloch, 1943).

As an essential part of the LingTester is the Frailty predictive model, the examination of the most common classifiers for text classification was conducted. The constructed dataset was used to feed the classifier using only the currently optimal features:

- transcript
- sex
- number\_of\_words
- text\_entropy
- desc\_event\_eng\_sentiment
- prev\_text\_eng\_sentiment
- desc\_image\_mispelled
- desc\_event\_mispelled
- class

For the model evaluation, the well known cross validation technique was deployed. Cross validation assesses how the results of a statistical analysis will generalize to an independent dataset. It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice. Due to the lack of sufficient examples (only 111 labeled instances in the dataset), the common 10-fold cross validation approach was dropped and the Leave-one-out cross-validation (LOOCV) was used instead. In LOOCV is a particular case of leave-p-out cross-validation with  $p = 1$ , where a statistic on the left-out samples is computed.

The next table summarises the accuracies obtained by the trained models, as can be seen in the table, the Decision Tree model scores the best results, thus it is selected as a first classification approach for the LingTester.

Classifier Name	Accuracy
Decision Tree (C4.5)	62.16 %
NaiveBayes	54.05 %
K-Nearest Neighbours	61.26 %
Support Vector Machines	31.35 %
Neural Networks	61.92 %

Table 3. Summary of classification algorithms

**Model optimization**

Before embedding the Decision Tree Classifier to the LingTester the final step of model parameter optimization was conducted using the Weka Data mining and Classification tool(Holmes, 1994). Specifically, the software enables the parameterization of the Decision Tree model using a series of eleven parameters, a table with the most important parameters and their description follows below.

Parameter Name	Description
Binary splits	Whether to use binary splits on nominal attributes when building the trees.
Confidence factor	The confidence factor used for pruning (smaller values incur more pruning).
MinNumObj	The minimum number of instances per leaf.
Reduced Error Pruning	Whether reduced-error pruning is used instead of C.4.5 pruning.
Unpruned	Whether pruning is performed.
Use Laplace	Whether counts at leaves are smoothed based on Laplace.

Table 4. Parameterization of the Decision Tree model

The process of model parameter optimization is a highly empirical process, although there have been some efforts in the field, for example Auto-Weka (Thornton, 2013). As this is a first approach of the classification task the simple strategy of test and recall has been followed. In order to improve the accuracy of LingTester the Train dataset was further investigated. In relation to its class (Non-frail, Pre-frail, Frail) it was judged as highly imbalanced as the Non-frail instances were only representing 12% of the class. For this reason a temporary decision was made to group the classes Non-frail and Pre-frail to restore the balance of the dataset until more

data is collected by the other tasks. After the overall model optimization an astonishing 10% accuracy increase was achieved. Figure 9 and Table 5 present the model statistics and the optimized parameter values accordingly.

```

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      80          72.0721 %
Incorrectly Classified Instances    31          27.9279 %
Kappa statistic                    0.3951
Mean absolute error                 0.3206
Root mean squared error             0.4467
Relative absolute error             66.3472 %
Root relative squared error         90.5213 %
Total Number of Instances          111

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
                0.836   0.455   0.737     0.836   0.783     0.742   prefrail
                0.545   0.164   0.686     0.545   0.608     0.742   frail
Weighted Avg.   0.721   0.339   0.717     0.721   0.714     0.742

=== Confusion Matrix ===

 a  b  <-- classified as
56 11 | a = prefrail
20 24 | b = frail
    
```

Figure 9. Model statistics

Parameter Name	Parameter Value
Binary splits	False
Confidence factor	0.25
MinNumObj	2
Reduced Error Pruning	False
Unpruned	True
Use Laplace	False

Table 5. Parameterization values of the Decision Tree model



## 6. Conclusions

The obtained accuracy of 72% by the common Decision Tree classifier seems promising and is a good starting point for the construction of more complex ensemble models. Although many of the state of the art techniques were implemented for the feature extraction process, they can not be exploited at the current phase of the project due to the small amount of collected data. As the dataset grows it is expected more features to contribute to the performance of the predictive model therefore a more capable and accurate model can be obtained and integrated to the FrailSafe user software package.

The further reduction of features and the hyper-optimization of the model showed to have slightly better test performance but in practice is an overfit (Everitt, 2002) of the current dataset and it will probably lead to worse overall results as the dataset examples increase in the future.



## 7. System development

### 7.1 Software used for training and prediction

Software development was based on various programming languages along with the use of third party tools, for both database management, creating models and predicting results. First steps of database management and export tools are based on Python (see Annex [12.1 Database management & feature extraction](#)), while the use of the prediction model is based on a Java implementation (see Annex [12.2 Prediction model](#)). Both programming languages are operating system independent, meaning that upon successful installation of Python console for the first, and Java Virtual machine for the later, provided programmes may start processing.

Furthermore, all third party tools within the discussed methodology should be provided. Firstly, the Greek POS tagger is a third party tool that needs Java Virtual Machine to run. It can be downloaded from the official site of Natural Language Processing Group at Department of Informatics - Athens University of Economics and Business (<http://nlp.cs.aueb.gr/software.html>). It is executed automatically from the Python script, whenever we request POS information for existing or new patient data.

For the classification process WEKA suite was used. Weka is a collection of machine learning algorithms for data mining tasks. The algorithms can either be applied directly to a dataset or called from your own Java code. Weka contains tools for data pre-processing, classification, regression, clustering, association rules, and visualization. It is also well-suited for developing new machine learning schemes. In the following section, we discuss steps to reproduce all installations steps.

### 7.2 Installation notes for training and prediction

Extra steps must be carried out so as for all submodules to work correctly, for the current state of LingTester. While the following submodules are platform independent, this doesn't mean that steps are also the same. Each platform may request different dependencies, but this analysis will assume Linux as the underlying operating system. These are the main parts:

- Python core and Python modules
- Java Virtual Machine
- POS Tagger
- OpenOffice dictionaries
- Weka

### 7.2.1 Python core and Python modules

Python is preinstalled in all Linux distros. However, we must install some extra modules. The following modules can be easily installed using the `pip` package, by issuing a command of type `sudo pip install {module}` for each module. These instructions are also described within the source code. The modules are:

1. **pyenchant**, for the translation service
2. **httplib2**, to request access to the MyMemory service through http protocol
3. **nlTK**, to export tf-idf features
4. **pattern**, to export sentiment features in english

### 7.2.2 Java Virtual Machine (JVM)

In order to download and install the Java Virtual Machine, user must navigate to the download page<sup>1</sup>, click Download and execute the downloaded file. Installation will automatically finalise. JVM is needed in order to execute both the ready made model and the POS tagger, in case it is needed for new data.

### 7.2.3 POS Tagger

POS tagger is available for download from the official page<sup>2</sup>. Assuming JVM works as expected, downloading and installing all files as stated in the readme file from the downloaded archive is all that is needed to have access to this service.

### 7.2.4 OpenOffice dictionaries

While `pyenchant` python module gives access to the speller of OpenOffice, we have to install the speller for each language, in case it is not already installed. In Linux, we can install new dictionaries by issuing in the command line the command `sudo apt-get install myspell-gr-el`. In case we have a different package manager, we issue the install command for the `myspell-gr-el` package.

### 7.2.5 Weka

We can easily download and install WEKA by navigating to the page official<sup>3</sup> and following the detailed steps in that page, according to our operating system.

---

<sup>1</sup> <https://www.java.com/en/>

<sup>2</sup> <http://nlp.cs.aueb.gr/software.html>

<sup>3</sup> <http://www.cs.waikato.ac.nz/ml/weka/downloading.html>

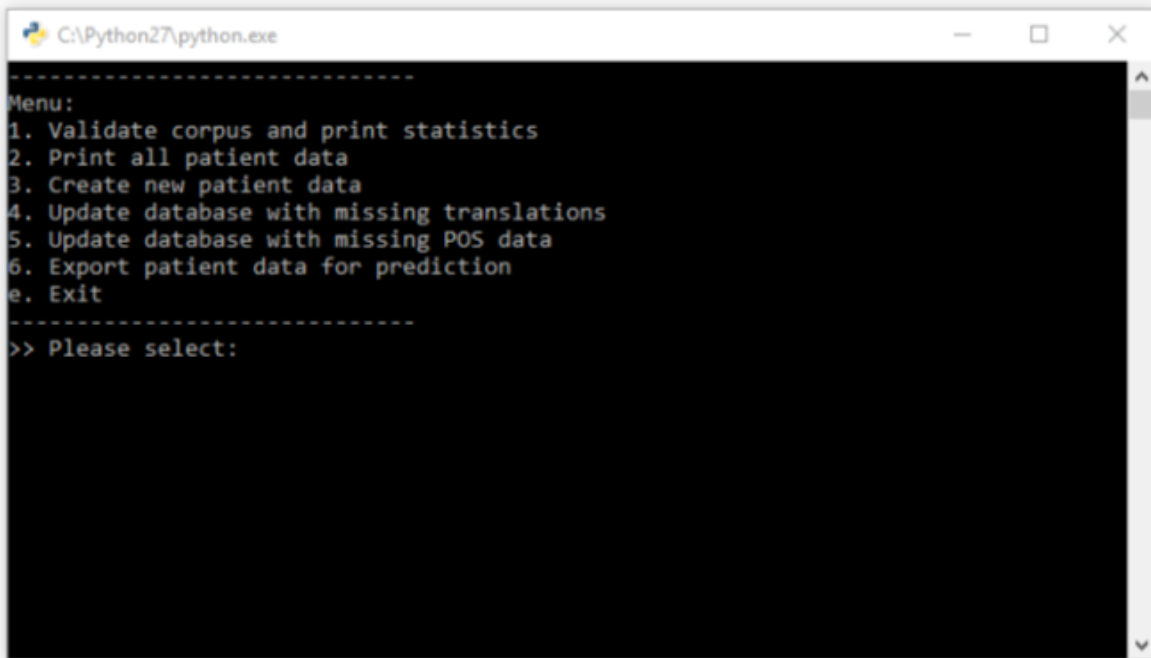
## 8. Frailsafe user software

### 8.1 Steps to import new data

There are two reasons to import new data within our offline database. First one is to populate new training data and reconstruct the prediction model, a process which cannot be done outside laboratory or by non technical persons. The reason for this, is that a full installation of WEKA is needed along with various parameters that make this procedure not feasible for everyday use. The second reason is to use the model that has already been trained and exported for easy use (see next chapter). For this to work, we execute the secondary file `offline-parser-ui.py` through the console, `python offline-parser-ui.py`, and the following menu appears:

1. *Validate corpus and print statistics*: This option goes through each patient data, and validates there are no missing or not accepted information
2. *Print all patient data*: This option requests a patient id from the user, and dumps all patient data within the console for easy access.
3. *Create new patient data*: This option will start by asking the ID of the new patient data and in case this ID already exists, then an error is shown to the user and he has to start over. This validation is needed in order to avoid overwriting existing data by accident.
4. *Update database with missing translations*: As the name suggests, this option will automatically update any missing translations
5. *Update database with missing POS data*: Also, as the name suggests, this will update whole database with any missing POS data for future use.
6. *Export patient data for prediction*: This option requests a patient id from the user and then creates the ARFF that is needed in order for the prediction model to work correctly (see next chapter).
7. *Exit*: This will force for the application to terminate

As this tool, is also Python based, it is cross-platform. However, this doesn't mean it can work out of the box. Extra steps need to be done in order for all submodule operate normally (See *Installation Notes*, under System development chapter).



```
C:\Python27\python.exe
-----
Menu:
1. Validate corpus and print statistics
2. Print all patient data
3. Create new patient data
4. Update database with missing translations
5. Update database with missing POS data
6. Export patient data for prediction
e. Exit
-----
>> Please select:
```

Figure 10. Screenshot of real time operation

## 8.2 Steps to export results

In order to for the FrailSafe user to obtain predictions for a number of subjects, a java software package was developed. The official name of the package is Predictor due to the task assigned to it.

The Predictor is a cross-platform command line tool that expects as input an arff file containing the pre-processed collected data of the subjects exported by the user using the Offline-parser-ui tool. Currently only the default ``in.arff'' can be used and Predictor expects it in the current working directory. An example input file structure is presented in the following figure.

```

relation frailsafe_111

@attribute transcript {yes,no,na}
@attribute class {prefrail,frail}
@attribute sex {male,female}
@attribute number_of_words numeric
@attribute text_entropy numeric
@attribute desc_event_ENG_sentiment numeric
@attribute prev_text_ENG_sentiment numeric
@attribute desc_image_misspelled numeric
@attribute desc_event_misspelled numeric

@data
na,?,female,48,4.023,0,0,0.429,0.571
na,?,female,43,3.869,0.625,0,0.345,0.467
na,?,male,0,0,0,0,0,0
na,?,female,33,3.805,1.25,0,0.4,0.368
na,?,male,0,0,0,0,0,0
na,?,female,0,0,0,0,0,0
    
```

Figure 11. Input file structure

In order for the Predictor not be computationally expensive, the training of the classifier at runtime of the tool was intentionally avoided. Instead, a pre-computed model of the classifier is delivered with the software package and is loaded by the tool.

The detailed program algorithm of the Predictor tool is presented below.

**Predictor tool algorithm**

**Input:**

- Load the pretrained model (M)
- Load the Test dataset (T)

Count the number of test instances (N)

**Loop for N**

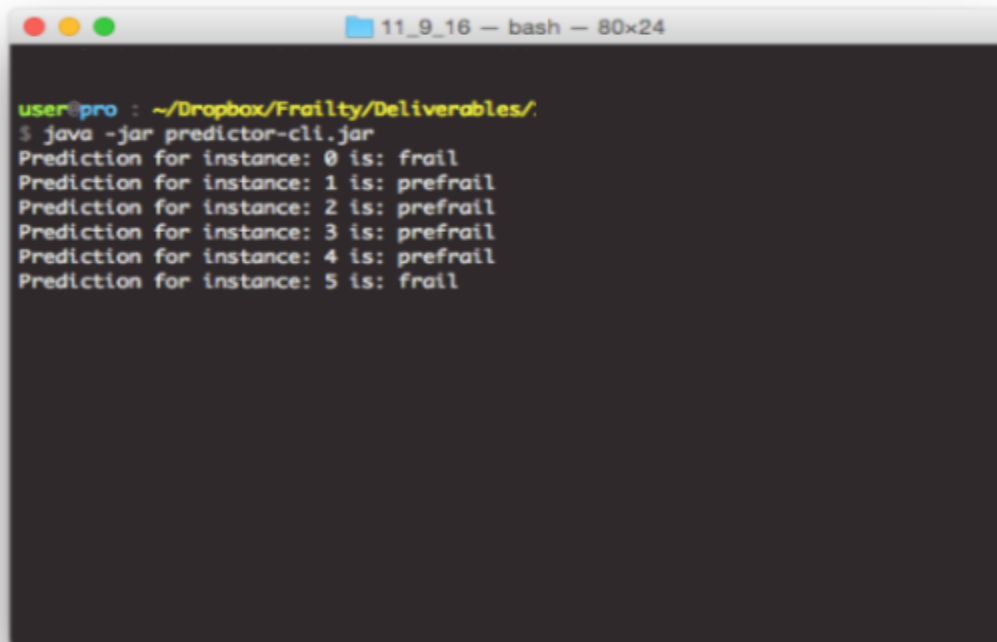
- Predict Instance T(i)
- Print prediction for T(i)

Software tool usage instructions:

1. Open a terminal inside the deliverables root directory
2. Run the command `java -jar predictor-cli.jar`

Notes: The user's system should use java version 1.6.0 or higher.

It is recommended to use the `-Xmx2g` java parameter if the input arff file has more than a few thousand instances.



```
11_9_16 — bash — 80x24
user@pro : ~/Dropbox/Frailty/Deliverables/
$ java -jar predictor-cli.jar
Prediction for instance: 0 is: frail
Prediction for instance: 1 is: prefrail
Prediction for instance: 2 is: prefrail
Prediction for instance: 3 is: prefrail
Prediction for instance: 4 is: prefrail
Prediction for instance: 5 is: frail
```

Figure 12. Screenshot of prediction programme while executing

## 9. Future improvements for the Prototype

A series of improvements is expected to improve the overall performance, add more predictive capabilities and improve the LingTester user experience. More specifically, on the subject of data collection, management and dataset exportation, a number of actions like the offline database population with new collected data, the features evaluation against a different input language (French language) and the optimization and bug fix of current features and feature extraction methodology is expected to improve the quality of the train dataset. On the subject of the classification task, the future research includes the exhaustive research on feature selection methodologies, the experimentation on supervised ensemble classifier models and the deployment of semi-supervised (Chapelle, 2006) techniques in order enhance the small available dataset with a bigger unlabeled (Ratsaby, 1995) dataset. Moreover, two new areas of research are going to be explored, the detection of suicidal statements and the patient mental state transition in time.

To benefit the end user of LingTester and simplify the prediction process, the integration of the two software tools, presented in section 8, in a complete and all inclusive software package will be implemented.

## 10. Ethics and Safety

Throughout this study's methodology, special care has been taken for ethical and safety issues. The nature of the study requires the processing, storing and analysis of a large amount of data. In all these stages, confidentiality and personal data protection will be reassured by an anonymization procedure. Each participant is traced solely by his/her ID, provided initially by the recruitment center, a number and only this, with no identifiable personal data, will be exposed to large scale data exchange like name, date of birth, place of living. Access to the database have only specific people, researchers, in order to create the prediction models.

The data persistence and analysis will comply with the data protection guidelines reported in deliverable "D9.9: Ethics, Safety and Health Barriers" (Section 6) with the aim of, at same time, keeping the maximum level of security and privacy of the data and allowing the successful performance of the other tasks of the project. Moreover, data will be obtained in accordance to the local ethics requirements. Any information regarding the participants will be treated as sensitive personal data (as defined in deliverable D9.9) and kept strictly private. Future provided data will be thoroughly checked by semi-automatic algorithms in order to anonymize any personal identifiers like full names, dates, emails, communication cellphone or landline numbers – hence falling outside the scope of legislation concerning personal data.



## 11. References

- A survey of text classification algorithms CC Aggarwal, CX Zhai - Mining text data, 2012 - Springer
- Gerard Salton, Edward A. Fox, and Harry Wu. 1983. Extended Boolean information retrieval. *Commun. ACM* 26, 11 (November 1983), 1022-1036. DOI=<http://dx.doi.org/10.1145/182.358466>
- Eibe, F., Mark, A. Hall, and Ian, H. Witten (2016). The WEKA Workbench. Online Appendix for "Data Mining: Practical Machine Learning Tools and Techniques", Morgan Kaufmann, Fourth Edition, 2016.
- K. Denecke, "Using SentiWordNet for multilingual sentiment analysis," Data Engineering Workshop, 2008. ICDEW 2008. IEEE 24th International Conference on, Cancun, 2008, pp. 507-512. doi: 10.1109/ICDEW.2008.4498370
- Koleli, E. A new Greek part-of-speech tagger, based on a maximum entropy classifier. Diss. Thesis, Athens University of Economics, 2011.
- Tsakalidis, A., Papadopoulos, S. and Kompatsiaris I. (2014) An Ensemble Model for Cross-Domain Polarity Classification on Twitter. Proceedings of 15th International Conference on Web Information Systems Engineering (WISE 2014), LNCS, Part II, pp. 168-177, Thessaloniki, Greece, October 12-14, 2014, Springer
- Gareth James; Daniela Witten; Trevor Hastie; Robert Tibshirani (2013). An Introduction to Statistical Learning. Springer.
- Bermingham, Mairead L.; Pong-Wong, Ricardo; Spiliopoulou, Athina; Hayward, Caroline; Rudan, Igor; Campbell, Harry; Wright, Alan F.; Wilson, James F.; Agakov, Felix; Navarro, Pau; Haley, Chris S. (2015). "Application of high-dimensional feature selection: evaluation for genomic prediction in man". *Sci. Rep.*
- Pearl, J. Heuristics: Intelligent Search Strategies for Computer Problem Solving. Addison-Wesley, 1984.
- Khachatryan, A.; Semenovskaya, S.; Vainshtein, B. (1979). "Statistical-Thermodynamic Approach to Determination of Structure Amplitude Phases". *Sov.Phys. Crystallography*.
- Mitchell, Melanie (1996). An Introduction to Genetic Algorithms. Cambridge, MA: MIT Press.
- Geisser, Seymour (1993). Predictive Inference. New York, NY: Chapman and Hall.
- Mitchell, Tom M. (1997). Machine Learning. The Mc-Graw-Hill Companies, Inc.
- Duda, Richard O.; Hart, Peter E.; Stork, David G. (2001). "Unsupervised Learning and Clustering". *Pattern classification* (2nd ed.). Wiley.
- Russell, Stuart; Norvig, Peter (2003) [1995]. Artificial Intelligence: A Modern Approach (2nd ed.). Prentice Hall.
- Cortes, C.; Vapnik, V. (1995). "Support-vector networks". *Machine Learning*.

- McCulloch, Warren; Walter Pitts (1943). "A Logical Calculus of Ideas Immanent in Nervous Activity". Bulletin of Mathematical Biophysics.
- G. Holmes; A. Donkin; I.H. Witten (1994). «Weka: A machine learning workbench». Proc Second Australia and New Zealand Conference on Intelligent Information Systems, Brisbane, Australia.
- Chris Thornton, Frank Hutter, Holger Hoos, and Kevin Leyton-Brown, Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms. In Proc. of KDD 2013
- Everitt B.S. (2002) Cambridge Dictionary of Statistics, CUP.
- Chapelle, Olivier; Schölkopf, Bernhard; Zien, Alexander (2006). Semi-supervised learning. Cambridge, Mass.: MIT Press.
- Ratsaby, J. and Venkatesh, S. Learning from a mixture of labeled and unlabeled examples with parametric side information. In Proceedings of the Eighth Annual Conference on Computational Learning Theory, pages 412-417 (1995).

## 12. Source files

These are the files that accompany this deliverable:

- Folder: demo
  - File: frailsafe.model
    - Pre-trained prediction model
  - File: in.arff
    - Test input data. This is a text file structured like [figure 11](#). For more details, please see chapter 8.1
  - File: predictor-cli.jar
    - This java software uses the pre-trained prediction model in order to predict patients' mental state from test data found in file in.arff(6 patients' test data after feature extraction).
  - File: README.txt
    - Readme file of how to execute the java file
- Folder: source code
  - File: PredictorCLI.java
    - Source code of the demo *predictor-cli.jar* file
  - File: offline-parser.py
  - File: offline-parser-ui.py
    - Source file in Python that manage the offline database while also support the feature extraction process

## 13. Annexes

### 13.1 Database management & feature extraction

The following are the contents of the base Python library with various functions that help utilise the offline database in its full extent.

```

1.  # -*- coding: utf-8 -*-
2.  """
3.
4.  @author: Charalampos Tsimpouris
5.  @author: Nikolaos Fazakis
6.
7.
8.  """
9.
10. # To install execute: pip install pyenchant
11. import enchant
12.
13. # To install execute: pip install httpplib2
14. import httpplib2
15. import json
16. import math
17.
18. # To install execute: pip install nltk
19. import nltk
20.
21. import os
22. import pickle
23.
24. # To install execute: pip install pattern
25. from pattern.en.wordnet import sentiment
26.
27. import stemming
28. import shutil
29. import subprocess
30.
31. # py-translate Cannot work
32. # it uses, free google services and is easilly blocked
33. # import translate
34.
35. # Google API
36. # cannot be used, only with billing plans
37. #from google.cloud import translate
38.
39. import warnings
40. from sklearn.feature_extraction.text import
    TfidfVectorizer
41.
42. # This is the path where all patient data is stored
43. # .. one file per patient, with file name from the
    patient id
44. data_path = './Data'
45.
46. frailsafe_google_api_key = '?????????'
47.
48. # MyMemory Translation
49. # .. these are basic settings for the translation
    service
50. mymemory_account_email = 'kifinas.uop@gmail.com'
51. mymemory_base_url =
    u'http://api.mymemory.translated.net/get'
52.
53. # These tags must always exist
54. # .. along with the available tags
55. # .. and they are class tags
56. verify_tags = {}
57. verify_tags['-transcript'] = ('yes', 'no', 'na')
58. verify_tags['-sex'] = ('male', 'female')
59. verify_tags['-tag'] = ('nonfrail', 'prefrail',
    'frail', 'na')
60. verify_tags['-language'] = ('greek', 'greek-
    polytonic', 'french', 'greek-cypriot')
61.
62. # Langue convert helper dictionaries for various
    reasons
63. langs = {}
64. langs['greek'] = 'el'
65. langs['greek-polytonic'] = 'el'
66. langs['french'] = 'fr'
67. langs['greek-cypriot'] = 'el-cy'
68.
69. langs_speller = {}
70. langs_speller['greek'] = 'el_GR'
71. langs_speller['greek-polytonic'] = 'el_GR'
72. langs_speller['french'] = 'fr'
73. langs_speller['greek-cypriot'] = 'el_GR'
74.
75. # These tags are multi line

```

```

76. # .. and make the multi text that we try to classify
    from
77. multi_line_tags = ['-desc_image', '-
    desc_event', '-prev_text']
78.
79. # Part of speech info, in case there is one
80. multi_line_tags_POS = ['-desc_image_POS', '-
    desc_event_POS', '-prev_text_POS']
81.
82. # English translation of the initial text
83. multi_line_tags_ENG = ['-desc_image_ENG', '-
    desc_event_ENG', '-prev_text_ENG']
84.
85. def split_list():
86.     """This function tries to split the main initial
    patient
87.     list in a simplistic way.
88.     Caution, detroyes existing data.
89.     """
90.     f = open(data_path + '/lists.txt', 'r')
91.     lines = f.readlines()
92.     f.close()
93.
94.     patients = {}
95.     for line in lines:
96.         if not line.startswith('-patient'):
97.             continue
98.
99.         tag, pid = line.strip().split(' ')
100.        if pid in patients:
101.            print "Patient %s already set"
% pid
102.            print "Aborting"
103.            return
104.            patients[ pid ] = line
105.
106.
107.            cpid = None
108.            pdata = []
109.
110.            for line in lines:
111.                line = line.strip()
112.
113.                if line.startswith('-patient'):
114.                    if cpid is not None:
115.                        f = open(data_path +
'/p.%s.txt' % cpid, 'w')
116.                        f.write("\n".join(pdata).strip() + "\n")
117.                        f.close()
118.
119.                        pdata = []
120.                        tag, cpid =
line.strip().split(' ')
121.
122.                        pdata.append(line)
123.
124.                        if cpid is not None:
125.                            f = open(data_path + '/p.%s.txt' %
cpid, 'w')
126.                            f.write("\n".join(pdata).strip() +
"\n")
127.                            f.close()
128.
129.                        def fetch_patient_ids():
130.                            """Identifies all patient ids, as
stated in the filenames
131.                            """
132.                            files = os.listdir(data_path)
133.                            out_files = []
134.                            for f in files:
135.                                if not f.startswith('p.') or not
f.endswith('.txt'):
136.                                    continue
137.
138.                                pre, pid, suf =
f.strip().split('.')
139.                                out_files.append(pid)
140.
141.                                out_files.sort()
142.                                return out_files
143.
144.                        def fetch_patient_data(cpid):
145.                            """Tries to load all patient data,
based on the id
146.                            """
147.                            try:
148.                                f = open(data_path + '/p.%s.txt' %
cpid, 'r')
149.                                lines = f.readlines()
150.                                f.close()
151.                                except:
152.                                    print 'Error opening patient data
file %s' % cpid
153.                                    return {}
154.
155.                                ret = {}
156.                                ctag = None
157.                                for line in lines:
158.                                    line = line.strip()
159.                                    if line.startswith('-'):
160.                                        if ctag:
161.                                            ret[ctag] =
ret[ctag].strip()
162.
163.                                            if line.find(' ') > 0:
164.                                                ctag = line[:line.find('
')]
.strip()
165.                                                info = line[line.find('
')+1:].strip()
166.                                            else:
167.                                                ctag = line

```

```

168.         info = ''
169.         ret[ctag] = ''
170.         ret[ctag] += info
171.         continue
172.
173.         ret[ctag] += line + "\n"
174.     if ctag:
175.         ret[ctag] = ret[ctag].strip()
176.
177.     return ret
178.
179.     def save_patient_data(cpid, cpdata):
180.         """Saves all patient data in a file
181.         """
182.         f = open(data_path + '/p.%s.txt' %
cpid, 'w')
183.
184.         cpdata['-patient'] = cpid
185.
186.         for tag in cpdata:
187.             if cpdata[tag] is None:
188.                 cpdata[tag] = ''
189.
190.             if tag in multi_line_tags or tag in
multi_line_tags_POS or tag in multi_line_tags_ENG:
191.                 f.write("%s\n\n" % tag)
192.                 f.write("%s\n\n" %
cpdata[tag].strip())
193.                 continue
194.
195.                 f.write("%s %s\n" % (tag,
str(cpdata[tag].strip())))
196.         f.close()
197.         return cpdata
198.
199.     def print_patient_data(pdata):
200.         """Tries to print all patient data, as
beautiful as it can
201.         """
202.         for k in pdata:
203.             print '%s:%s' % (k, pdata[k])
204.
205.     def validate_patient_data():
206.         """Validates there are no missing tags
207.         in all patient files, around the
required ones
208.         Also, it tries to print some minor
statistics
209.         """
210.         pids = fetch_patient_ids()
211.
212.         stats = {}
213.         for cpid in pids:
214.             cpdata = fetch_patient_data(cpid)
215.
216.             for t in verify_tags:
217.                 valid = verify_tags[t]
218.                 if not t in cpdata or cpdata[t]
is None or cpdata[t] == '':
219.                     print 'Patient %s is
missing %s data' % (cpid, t)
220.                     continue
221.
222.                     if cpdata[t] not in valid:
223.                         print 'Patient %s has
invalid %s data %s' % (cpid, t, cpdata[t])
224.                         continue
225.
226.                     if not t in stats:
227.                         stats[t] = {}
228.                         stats[t][ cpdata[t] ] =
stats[t].get(cpdata[t], 0) + 1
229.
230.                     print 'Checked %d patients, and here
are the statistics' % len(pids)
231.                     for t in verify_tags:
232.                         valid = verify_tags[t]
233.
234.                         print t
235.                         for k in valid:
236.                             print ' %s:%d' % (k,
stats[t].get(k, 0))
237.
238.         def temp_change_tags():
239.             """Makes a change of tags, after
initial estimation
240.             Temporary function
241.             """
242.             pids = fetch_patient_ids()
243.
244.             for cpid in pids:
245.                 cpdata = fetch_patient_data(cpid)
246.                 if '-desc_image' not in cpdata:
247.                     cpdata['-desc_image'] =
cpdata.get('-perigrافي_eikonas', '')
248.                     del cpdata['-
perigrافي_eikonas']
249.
250.                     if '-desc_event' not in cpdata:
251.                         cpdata['-desc_event'] =
cpdata.get('-perigrافي_gegonotos', '')
252.                         del cpdata['-
perigrافي_gegonotos']
253.
254.                     save_patient_data(cpid, cpdata)
255.
256.         def clean_up_text(text, lang = 'english'):
257.             """Tries to clean up text, as much as
possible

```

```

258.         in some ways language
           inndipandantly
259.         in some ways not
260.         """
261.         new_list = []
262.
263.         text = text.replace('.', '. ')
264.         # Minor clean up per language
265.         if lang.startswith('greek'):
266.             text = clean_greek_letters(text)
267.         elif lang == 'french':
268.             text = clean_french_letters(text)
269.
270.         text = upper(text)
271.
272.         words = get_words(text)
273.         for w in words:
274.             if len(w) <= 3:
275.                 continue
276.
277.             w = stemming.stem(w)
278.
279.             new_list.append(w)
280.
281.         return ' '.join(new_list)
282.
283.     def create_arff(relation = 'fraildata'):
284.         """Create arff for WEKA with all
           features available
285.         """
286.         pids = fetch_patient_ids()
287.
288.         out = []
289.         out.append('@RELATION %s' % relation)
290.         out.append('')
291.
292.         basic_tags = []
293.         for t in verify_tags:
294.             basic_tags.append(t)
295.             valid = verify_tags[t]
296.             tag = t.lstrip('-')
297.             if tag == 'tag':
298.                 tag = 'class'
299.                 valid = ('nonfrail',
                    'prefrail', 'frail')
300.             out.append('@ATTRIBUTE %s %s' %
                (tag, ','.join(valid)))
301.
302.         other_attributes = []
303.
304.         other_attributes.append('get_feature_length')
305.
306.         other_attributes.append('get_feature_number_of_sentences')
307.
308.         for attr in other_attributes:
309.             out.append('@ATTRIBUTE %s %s' %
                (globals()[attr](' ', 'title'), globals()[attr](' ',
                    'type')))
310.
311.         for tag in multi_line_tags_ENG:
312.             out.append('@ATTRIBUTE %s %s' %
                (tag.lstrip('-') + '_sentiment', 'real'))
313.
314.         for tag in multi_line_tags:
315.             out.append('@ATTRIBUTE %s %s' %
                (tag.lstrip('-') + '_misspelled', 'real'))
316.
317.         corpus = get_corpus()
318.
319.         labels = []
320.         texts = []
321.         text_POS = []
322.         for cpid in corpus:
323.             labels.append(cpid)
324.             texts.append(corpus[cpid]['clean_text'])
325.
326.             pos_data =
                 pos_explode_data(corpus[cpid]['text_POS'])
327.             temp_pos_as_text = []
328.             # ta panta ola
329.             for word, ps in pos_data:
330.                 temp_pos_as_text.append('8'.join(ps))
331.
332.             # 1:1 mono ta prwta
333.             for word, ps in pos_data:
334.                 temp_pos_as_text.append(ps[0])
335.
336.             # 1:1 mono ta prwta/deftera
337.             for word, ps in pos_data:
338.                 temp_pos_as_text.append('8'.join(ps[0:1]))
339.
340.             # 1:1 mono ta prwta/deftera/trita
341.             for word, ps in pos_data:
342.                 temp_pos_as_text.append('8'.join(ps[0:2]))
343.
344.             # 1:1 mono ta deftera
345.             for word, ps in pos_data:
346.                 temp_pos_as_text.append(ps[0])

```

```

347.
348.         # 1:1 mono ta trita
349.         for word, ps in pos_data:
350.             temp_pos_as_text.append(ps[0])
351.
352.             text_POS.append('
'.join(temp_pos_as_text).replace('-', ''))
353.
354.         # TF-IDF on POS data
355.         tf_POS =
TfidfVectorizer(analyzer='word', ngram_range=(1,2),
min_df = 0)
356.         tfidf_matrix_POS =
tf_POS.fit_transform(text_POS)
357.         feature_names_POS =
tf_POS.get_feature_names()
358.
359.         for i in range(len(feature_names_POS)):
360.             out.append('@ATTRIBUTE tf-pos-%d
real %% %s' % (i, 'POS: ' + feature_names_POS[i]))
361.             dense_POS = tfidf_matrix_POS.todense()
362.
363.         # TF-IDF on stemmed text
364.         tf = TfidfVectorizer(analyzer='word',
ngram_range=(1,2), min_df = 0)
365.         tfidf_matrix = tf.fit_transform(texts)
366.         feature_names = tf.get_feature_names()
367.
368.         for i in range(len(feature_names)):
369.             out.append('@ATTRIBUTE tf-%d real
%% %s' % (i, feature_names[i]))
370.             dense = tfidf_matrix.todense()
371.
372.             out.append('')
373.             out.append('@DATA')
374.
375.             for i in range(len(labels)):
376.                 cpid = labels[i]
377.
378.                 # Patients with inknown frailty tag
are automatically removed
379.                 if not '-tag' in
corpus[cpid]['data'] or corpus[cpid]['data']['-tag']
== 'na':
380.                     print "Skipping patient with ID
%s: missing frailty tag" % cpid
381.                     continue
382.
383.                 row = []
384.                 for bt in basic_tags:
385.                     row.append(
corpus[cpid]['data'].get(bt, 'na') )
386.
387.                 for attr in other_attributes:
388.                     clang = corpus[cpid]['data']['-
language']
389.
390.                     # To absorb all Greek
variations
391.                     if clang.startswith('greek'):
392.                         clang = 'greek'
393.
394.                     row.append(str(globals()[attr](corpus[cpid]['text'],
lang = clang)))
395.
396.                     # Sentiment score is based in the
english translation
397.                     for tag in multi_line_tags_ENG:
398.                         row.append(get_feature_sentiment_score(corpus[cpid]['d
ata'].get(tag, '')))
399.
400.                     for tag in multi_line_tags:
401.                         row.append(get_feature_mispelling_score(corpus[cpid]['
data'].get(tag, ''), lang = clang))
402.
403.                     # tf-idf info based on POS data
404.                     p_POS = dense_POS[i].tolist()[0]
405.                     for fi in
range(len(feature_names_POS)):
406.                         row.append('%0.3f' % p_POS[fi])
407.
408.                     # tf-idf based on stemmed data
409.                     p = dense[i].tolist()[0]
410.                     for fi in
range(len(feature_names)):
411.                         row.append('%0.3f' % p[fi])
412.
413.                     out.append(', '.join(row))
414.
415.                     filename = 'ARFFS/%s.arff' % relation
416.                     f = open(filename, 'w')
417.                     f.write("\n".join(out).encode('utf8'))
418.                     f.close()
419.
420.                     print 'Saved at %s.' % filename
421.
422.                 def get_corpus():
423.                     """Returns all corpus in a more
scientific-friendly way
424.                     """
425.                     pids = fetch_patient_ids()
426.
427.                     corpus = {}
428.
429.                     for cpid in pids:
430.                         cpdata = fetch_patient_data(cpid)
431.                         corpus[cpid] = {}
432.                         corpus[cpid]['data'] = cpdata

```



```

433.             corpus[cpid]['tag'] = cpdata['-
tag']
434.
435.             corpus[cpid]['text'] = ''
436.             for m in multi_line_tags:
437.                 corpus[cpid]['text'] +=
cpdata.get(m, ' ')
438.             corpus[cpid]['text'] =
corpus[cpid]['text'].strip()
439.
440.             corpus[cpid]['text_POS'] = ''
441.             for m in multi_line_tags_POS:
442.                 corpus[cpid]['text_POS'] +=
"\n" + cpdata.get(m, ' ')
443.
444.             corpus[cpid]['text_ENG'] = ''
445.             for m in multi_line_tags_ENG:
446.                 corpus[cpid]['text_ENG'] +=
"\n" + cpdata.get(m, ' ')
447.
448.             clang = corpus[cpid]['data']['-
language']
449.
450.             # To absorb all Greek variations
451.             if clang.startswith('greek'):
452.                 clang = 'greek'
453.
454.             tutf8 =
corpus[cpid]['text'].decode('utf-8')
455.             corpus[cpid]['clean_text'] =
clean_up_text(tutf8, clang)
456.
457.             return corpus
458.
459.             def get_feature_length(text, meta = None,
lang = 'english'):
460.                 if meta == 'title':
461.                     return 'text_length'
462.                 if meta == 'type':
463.                     return 'integer'
464.
465.                 return len(text)
466.
467.             def get_sentences(text, lang = 'english'):
468.                 """This should be language dependant to
be more precise
469.                 """
470.                 sent_tok_file =
"greek.law.utf8.70.pickle"
471.
472.                 f = open(sent_tok_file)
473.                 sent_tokenizer = pickle.load(f)
474.                 f.close()
475.
476.                 return
sent_tokenizer.tokenize(text.decode('utf8'))
477.
478.                 def get_feature_number_of_sentences(text,
meta = False, lang = 'english'):
479.                     """Returns a feature with number of
sentences
480.                     TODO: This has to be more elaborate
481.                     """
482.                     if meta == 'title':
483.                         return 'number_of_sentences'
484.                     if meta == 'type':
485.                         return 'integer'
486.
487.                     return len(get_sentences(text, lang))
488.
489.                 def get_words(text, lang = 'english'):
490.                     """Splits text in words
491.                     """
492.                     word_tokenizer =
nltk.WhitespaceTokenizer()
493.
494.                     # Word tokenizer, auto parses sentences
495.                     # ..so no need to split in sentences
496.                     return word_tokenizer.tokenize(text)
497.
498.                 def get_feature_word_count(text, meta =
False, lang = 'english'):
499.                     """Returns a feature with number of
words
500.                     TODO: This has to be more elaborate
501.                     """
502.                     if meta == 'title':
503.                         return 'number_of_words'
504.                     if meta == 'type':
505.                         return 'integer'
506.
507.                     return len(get_words(text, lang))
508.
509.                 def get_feature_words_per_sentence(text,
meta = False, lang = 'english'):
510.                     """Returns a feature with number of
words
511.                     TODO: This has to be more elaborate
512.                     """
513.                     if meta == 'title':
514.                         return
'number_of_words_per_sentence'
515.                     if meta == 'type':
516.                         return 'real'
517.
518.                     if int(get_feature_length(text)) <= 0:
519.                         return 0
520.

```

```

521.         return '%.3f' %
           (float(get_feature_word_count(text)) /
            get_feature_number_of_sentences(text))
522.
523.     def get_feature_text_shannon_entropy(text,
           meta = False, lang = 'english'):
524.         """Returns bits of entropy represented
           in a given string, per
525.         http://en.wikipedia.org/wiki/Entropy_(information_theo
           ry)
526.         """
527.         if meta == 'title':
528.             return 'text_entropy'
529.         if meta == 'type':
530.             return 'real'
531.
532.         mmap = {}
533.         for c in text:
534.             mmap[c] = mmap.get(c, 0) + 1
535.
536.         text_len = get_feature_length(text)
537.         result = 0.0
538.
539.         for c in mmap:
540.             freq = mmap[c] / float(text_len)
541.             result -= freq * (math.log(freq) /
           math.log(2))
542.
543.         return '%.3f' % result
544.
545.     def get_feature_sentiment_score(text, meta
           = False, lang = 'english'):
546.         """Returns sentiment score, works based
           on the english translation
547.         """
548.         if meta == 'title':
549.             return 'sentiment_score'
550.         if meta == 'type':
551.             return 'real'
552.
553.         v = 0
554.         for w in text.split(" "):
555.             w =
           w.strip(",.!?)#;:\'\"").lower()
556.             if w in sentiment:
557.                 v = v + sentiment[w][0] -
           sentiment[w][1]
558.         return str(v)
559.
560.     def get_feature_mispelling_score(text, meta
           = False, lang = 'english'):
561.         """Returns mispelling statistics
562.         """
563.         if meta == 'title':
564.             return 'mispelling_score'
565.         if meta == 'type':
566.             return 'real'
567.
568.         if not lang in langs_speller:
569.             warnings.warn('Unknown input
           language: %s' % (from_lang))
570.             return ''
571.             slang = langs_speller[lang]
572.
573.             word_counting = 0
574.             misspelled_words = 0
575.
576.             d = enchant.Dict(slang)
577.             for w in get_words(text, lang):
578.                 word_counting += 1
579.                 if not d.check(w):
580.                     misspelled_words += 1
581.
582.             if word_counting <= 0:
583.                 return '0.0'
584.
585.             return '%.3f' %
           (float(misspelled_words) / float(word_counting))
586.
587.     def clean_greek_letters(text):
588.         text = text.replace(u'ϊ', u'ι')
589.         text = text.replace(u'ϋ', u'υ')
590.         text = text.replace(u'Ϙ', u'χ')
591.         text = text.replace(u'ϙ', u'ψ')
592.         text = text.replace(u'Ϛ', u'φ')
593.         text = text.replace(u'ϛ', u'φ')
594.         text = text.replace(u'Ϝ', u'φ')
595.         text = text.replace(u'ϝ', u'φ')
596.         text = text.replace(u'Ϟ', u'φ')
597.         text = text.replace(u'ϟ', u'φ')
598.         text = text.replace(u'Ϡ', u'φ')
599.         text = text.replace(u'ϡ', u'φ')
600.         text = text.replace(u'Ϣ', u'φ')
601.         text = text.replace(u'ϣ', u'φ')
602.         text = text.replace(u'Ϥ', u'φ')
603.         text = text.replace(u'ϥ', u'φ')
604.         text = text.replace(u'Ϧ', u'φ')
605.         text = text.replace(u'ϧ', u'φ')
606.         text = text.replace(u'Ϩ', u'φ')
607.         text = text.replace(u'ϩ', u'φ')
608.         text = text.replace(u'Ϫ', u'φ')
609.         text = text.replace(u'ϫ', u'φ')
610.         text = text.replace(u'Ϭ', u'φ')
611.         text = text.replace(u'ϭ', u'φ')
612.         text = text.replace(u'Ϯ', u'φ')
613.         text = text.replace(u'ϯ', u'φ')
614.         text = text.replace(u'ϰ', u'φ')
615.         text = text.replace(u'ϱ', u'φ')
616.         text = text.replace(u'ϲ', u'φ')
617.         text = text.replace(u'ϳ', u'φ')
618.         text = text.replace(u'ϴ', u'φ')

```

```

619.         text = text.replace(u'İŽ', u'İ..')
620.         text = text.replace(u'İ«', u'İ..')
621.         text = text.replace(u'İ@', u'İ&')
622.
623.         text = text.replace(u'İ-', u'İ±')
624.         text = text.replace(u'İ-', u'İµ')
625.         text = text.replace(u'İ@', u'İ..')
626.         text = text.replace(u'İ', u'İ+')
627.         text = text.replace(u'İ@', u'İç')
628.         text = text.replace(u'İ@', u'İ.')
629.         text = text.replace(u'İŽ', u'İ&')
630.         text = text.replace(u'İS', u'İ±')
631.         text = text.replace(u'İ<', u'İ..')
632.         text = text.replace(u'İ@', u'İ+')
633.         text = text.replace(u'İ°', u'İ..')
634.
635.         text = text + u' '
636.         # text = text.replace(u'İ% ', u' ') #
        p.x. to "İµİçİçİ@İİİ±İİ" ginetai "İµİçİçİ@İİİ±İ±"
637.         text = text.replace(u'İf ', u'İ, ') #
638.         text = text.strip()
639.
640.         return text
641.
642.     def clean_french_letters(text):
643.         """Currently not implemented"""
644.         return text
645.
646.
647.     def upper(text):
648.         """Capitalize text"""
649.         text = text.replace(u'İ±', u'İ+')
650.         text = text.replace(u'İ±', u'İ+')
651.         text = text.replace(u'İ³', u'İ³')
652.         text = text.replace(u'İ', u'İ+')
653.         text = text.replace(u'İµ', u'İ.')
654.         text = text.replace(u'İç', u'İ-')
655.         text = text.replace(u'İ.', u'İ-')
656.         text = text.replace(u'İ.', u'İ+')
657.         text = text.replace(u'İ±', u'İ³')
658.         text = text.replace(u'İ@', u'İS')
659.         text = text.replace(u'İ»', u'İ.')
660.         text = text.replace(u'İŽ', u'İ@')
661.         text = text.replace(u'İŽ', u'İ@')
662.         text = text.replace(u'İç', u'İç')
663.         text = text.replace(u'İç', u'İç')
664.         text = text.replace(u'İ@', u'İ ')
665.         text = text.replace(u'İ@', u'İ±')
666.         text = text.replace(u'İf', u'İ@')
667.         text = text.replace(u'İ.', u'İ+')
668.         text = text.replace(u'İµ', u'İµ')
669.         text = text.replace(u'İ..', u'İç')
670.         text = text.replace(u'İ+', u'İ+')
671.         text = text.replace(u'İ+', u'İç')
672.         text = text.replace(u'İ', u'İ+')

```

```

673.         text = text.replace(u'İ&', u'İ@')
674.
675.         return text
676.
677.         def get_pos_info(text, debug = '',
        pos_directory = '/media/xaris/Data/PhD/POS/bin'):
678.             """Tries to execute POS tagger, and
        retrieves the results
679.             from the exported file"""
680.
681.             if not text or text == '':
682.                 return ''
683.
684.             # Till a better solution
685.             # .. everything has to be done where
        the java files are
686.             # .. so we hardcode pos_directory
687.             # .. but keep it as a parameter
688.
689.             current_directory = os.getcwd()
690.             os.chdir(pos_directory)
691.
692.             # this is the file where data will be
        stored
693.             in_file = pos_directory + '/in.txt'
694.
695.             f = open(in_file, 'w')
696.             f.write(text)
697.             f.close()
698.
699.             # output file is hardcoded according to
        maintainer
700.             out_file = pos_directory +
        '/result.txt'
701.
702.             # Keep everything clean
703.             f = open(out_file, 'w')
704.             f.write('')
705.             f.close()
706.
707.             res = subprocess.call(['java', '-jar',
        'POStagger.jar', '1', in_file])
708.             if res != 0:
709.                 # This means that program
        terminated with error
710.                 return ''
711.
712.             if not os.path.exists(out_file):
713.                 # I don't why this can happen
714.                 warnings.warn('Output file from POS
        method was empty, debug data: %s' % (str(debug)))
715.                 return ''
716.
717.             f = open(out_file, 'r')

```

```

718.         ret = f.readlines()
719.         f.close()
720.
721.         # And resoter current working directory
722.         os.chdir(current_directory)
723.
724.         return "".join(ret)
725.
726.     def get_translated_data(data, from_lang,
727.        to_lang = 'en', debug = ''):
728.         """Tries to translate the text to
729.         english
730.         using Google Translate API
731.         """
732.         if not from_lang in langs:
733.             warnings.warn('Unknown input
734.             language: %s' % (from_lang))
735.             return ''
736.             from_lang = langs[from_lang]
737.
738.             if from_lang == to_lang:
739.                 # Apparently there is no need to
740.                 call the API
741.                 return data
742.
743.                 # The following is based on py-
744.                 translate
745.                 # and is blocked for overuse
746.                 # -----
747.                 # return
748.                 translate.translator(from_lang, to_lang, data)
749.                 # -----
750.
751.                 # The following is based on Google API
752.                 # and is only on paid services
753.                 # -----
754.                 # translate_client =
755.                 translate.Client(frailsafe_google_api_key)
756.                 # translation =
757.                 translate_client.translate(data, source_language =
758.                 from_lang, target_language = to_lang)
759.                 # print('Text: {}'.format(text))
760.                 # print('Translation:
761.                 {}'.format(translation['translatedText'].encode('utf-
762.                 8')))
763.                 # -----
764.
765.                 # The following is based at MyMemory
766.                 service
767.                 # -----
768.                 lines = data.split("\n")
769.                 trans_result = ''
770.                 for line in lines:

```

```

759.             line = line.strip()
760.             if line == '':
761.                 continue
762.
763.             f = {}
764.             f['q'] = line
765.             f['langpair'] = '%s%s' %
766.             (from_lang, to_lang)
767.             f['of'] = 'json'
768.             f['de'] = mymemory_account_email
769.
770.             resp, json_content =
771.             httpLib2.Http().request("%s%s" % (mymemory_base_url,
772.             urllib.urlencode(f)))
773.             result = json.loads(json_content)
774.
775.             if result['responseStatus'] != 200:
776.                 print "Error from mymory,
777.                 aborting.. Debug: %s" % debug
778.                 return ''
779.
780.                 trans_result +=
781.                 result['responseData']['translatedText'] + "\n"
782.                 # -----
783.
784.                 return trans_result
785.
786.             def
787.             update_corpus_with_translations(force_rebuild =
788.             False):
789.                 """Get all text data from all patients
790.                 and updates the corpus with missing
791.                 translations
792.                 In order to avoid overuse of the
793.                 third-part service,
794.                 we save locally the translation for
795.                 future use.
796.                 The force_rebuild parameter, will
797.                 force to update all translations
798.                 """
799.                 pids = fetch_patient_ids()
800.
801.                 for cpid in pids:
802.                     cpdata = fetch_patient_data(cpid)
803.                     updated = False
804.                     for mt in multi_line_tags:
805.                         d = cpdata.get(mt, '')
806.                         # Adeio keimeno
807.                         if d == '' and not
808.                         force_rebuild:
809.                             continue
810.
811.                             # Exw idi ipologisei POS data
812.                             if cpdata.get(mt + '_ENG', '')
813.                             != '' and not force_rebuild:

```

```

802.         continue
803.
804.         # Den to exoume, as to paroume
805.         ret = get_translated_data(d,
cpdata['-language'], debug = cpid)
806.         if ret == '' and not
force_rebuild:
807.             continue
808.
809.             updated = True
810.             cpdata[mt + '_ENG'] = ret
811.
812.             # Something changed, time to store
it
813.             if updated:
814.                 print 'Updating patient %s' %
cpid,
815.                 save_patient_data(cpid, cpdata)
816.                 print '..Done'
817.
818.         def get_sentiment_analysis_greek(data, lang
= 'greek'):
819.             """The following idea was based on a
locally saved, sentiment analysis file
820.             .. but due to poor results, this
idea got rejected
821.             .. and we shifted to the sentiment
analysis of the english translation
822.             """
823.             # Load all data from sentiment corpus
824.             f = open('%s-sentiment-
lexicon/%s_sentiment_lexicon.tsv' % (lang, lang), 'r')
825.             lines = f.readlines()
826.             f.close()
827.
828.             titles = lines[0].lower().split('\t')
829.             del lines[0]
830.
831.             sentiment_words = {}
832.             for line in lines:
833.                 parts = line.split('\t')
834.
835.                 parts[0].rtim('-î -îç')
836.                 term = stemming.stem(parts[0])
837.
838.                 sentiment_words[ term ] = {}
839.                 for p in range(len(parts)):
840.                     sentiment_words[ term ][
titles[p] ] = parts[p]
841.
842.                 for k in sentiment_words.keys():
843.                     print k,
844.                     print
845.
846.                 data_clean = clean_up_text(data, lang)

```

```

847.         result = {}
848.         result['anger'] = 0
849.         result['disgust'] = 0
850.         result['fear'] = 0
851.         result['happiness'] = 0
852.         result['sadness'] = 0
853.         result['surprise'] = 0
854.
855.         words_identified = 0
856.
857.         for word in data_clean.split(' '):
858.             if not word in sentiment_words:
859.                 continue
860.
861.                 words_identified += 1
862.                 for r in result:
863.                     for k in sentiment_words[word]:
864.                         if not k.startswith(r) or
sentiment_words[word][k] == 'N/A':
865.                             continue
866.
867.                             result[r] +=
sentiment_words[word][k]
868.
869.                             if (words_identified > 0):
870.                                 for r in result:
871.                                     result[r] = float(result[r]) /
float(words_identified)
872.
873.                                     return result
874.
875.                 def
update_pos_info_everywhere(force_rebuild = False):
876.                     """Get all text data from all patients
and updates the corpus with Part-Of-
Speech information
877.
878.                     All results are saved within the
database
879.                     """
880.                     pids = fetch_patient_ids()
881.
882.                     for cpid in pids:
883.                         cpdata = fetch_patient_data(cpid)
884.                         updated = False
885.                         for mt in multi_line_tags:
886.                             d = cpdata.get(mt, '')
887.                             # Text is empty
888.                             if d == '' and not
force_rebuild:
889.                                 continue
890.
891.                                 # I already have this POS data
892.                                 if cpdata.get(mt + '_POS', '')
!= '' and not force_rebuild:

```

```

894.         continue
895.
896.         # POS data is missing, let's
calculate it
897.         ret = get_pos_info(d)
898.         if ret == '' and not
force_rebuild:
899.             continue
900.
901.         updated = True
902.         cpdata[mt + '_POS'] = ret
903.
904.         # Something changed, time to store
it
905.         if updated:
906.             print 'Updating patient %s' %
cpid,
907.             save_patient_data(cpid, cpdata)
908.             print '..Done'
909.
910.         def pos_explode_data(data):
911.             """Explodes all POS data from a string
912.             .. as given by the POS tagger
913.             .. and returns a more programming-
friendly object.
914.
915.             In case POS tagger changes, this
function must re-implemented
916.             """
917.             result = []
918.             lines = data.split("\n")
919.             for l in lines:
920.                 if l.strip() == '' or l.find(' ') <
0:
921.                     continue
922.
923.                 word, tags = l.split(' ')
924.                 ps = tags.split('/')
925.                 result.append((word, ps))
926.
927.             return result
928.
929.         def print_all_possible_pos_tags():

```

```

930.         """Prints all POS data within our
corpus
931.         for statistical reasons"""
932.         pids = fetch_patient_ids()
933.
934.         per_place = {}
935.         for cpid in pids:
936.             cpdata = fetch_patient_data(cpid)
937.             for mt in multi_line_tags_POS:
938.                 d = cpdata.get(mt, '')
939.                 # Adeio keimeno
940.                 if d == '':
941.                     continue
942.
943.                 pos_data = pos_explode_data(d)
944.                 for word, ps in pos_data:
945.                     for pos in range(len(ps)):
946.                         info = ps[pos]
947.                         if not pos in
per_place:
948.                             per_place[pos] = {}
949.
950.                         per_place[pos][ info ]
= per_place[pos].get(info, 0) + 1
951.
952.                         i = 0
953.                         while i in per_place:
954.                             print 'Position %d' % i
955.                             keys = per_place[i].keys()
956.                             keys.sort()
957.                             for k in keys:
958.                                 print ' ' * 3 + '%s: %d' % (k,
per_place[i][k])
959.
960.                                 i += 1
961.
962.         def install_spell_greek_checker_files():
963.             """This must be run a root
964.             """
965.             # Linux
966.             # sudo apt-get install myspell-gr-el
967.             pass

```

## 13.2 Prediction tool

```
1. package predictor;
2.
3. import weka.classifiers.Classifier;
4. import weka.core.Instances;
5. import weka.core.converters.ConverterUtils.DataSource;
6.
7. public class PredictorCLI {
8.
9.     public static void main(String[] args) {
10.
11.         Classifier cls;
12.         try {
13.             //load model
14.             cls = (Classifier) weka.core.SerializationHelper.read("frailsafe.model");
15.
16.
17.             DataSource source;
18.             try {
19.                 //load test data
20.                 source = new DataSource("in.arff");
21.                 Instances data = source.getDataSet();
22.                 if (data.classIndex() == -1)
23.                     data.setClassIndex(1); //class attribute is the second attribute
24.
25.                 //predict & print
26.                 for(int i=0; i<data.numInstances();i++){
27.                     double value=cls.classifyInstance(data.instance(i));
28.                     String prediction=data.classAttribute().value((int)value);
29.                     System.out.println("Prediction for instance: "+i+" is: "+prediction);
30.                 }
31.
32.             } catch (Exception e) {
33.                 // TODO Auto-generated catch block
34.                 e.printStackTrace();
35.             }
36.         } catch (Exception e) {
37.             // TODO Auto-generated catch block
38.             e.printStackTrace();
39.         }
40.     }
41.
42. }
```