



MASTER OF SCIENCE
IN ENGINEERING

Hes·SO

Haute Ecole Spécialisée
de Suisse occidentale

Fachhochschule Westschweiz

University of Applied Sciences and Arts
Western Switzerland

Master of Science HES-SO in Engineering
Av. de Provence 6
CH-1007 Lausanne

Master of Science HES-SO in Engineering

Orientation: Information and Communication Technologies (ICT)

GenBot

Author:

Romain Claret

`romain.claret@master.hes-so.ch`

Under the direction of:

Prof. Dr. Jean Hennebert

HES-SO//Fribourg

Institute of Complex Systems (iCoSys)

External expert:

Dr. Christophe Gisler

Lausanne, HES-SO//Master, May 27, 2019

Accepted by the HES-SO//Master (Switzerland, Lausanne) on a proposal from:

Prof. Dr. Jean Hennebert, deepening project supervisor
Dr. Christophe Gisler, iCoSys, main expert

Place, date: _____

Prof. Dr. Jean Hennebert
Supervisor

Prof. Dr. Philippe Passeraub
Dean, HES-SO//Master

WHO IS THE MR

Contents

Contents	v
Acknowledgements	ix
Acronyms	xi
Abstract	xiii
1 Introduction	1
1.1 Aim of Study	1
1.2 Scope and Study Borders	1
1.3 Industrial Interest	1
2 Questions	3
2.1 Initial and Broad Questions	3
2.2 Narrowed Questions	4
2.3 Potential Red lines	4
2.4 The Deepening Project Question and Red line	4
3 Plan	5
3.1 Constraints	5
3.2 Initial Plan	5
3.2.1 Tasks	5
3.2.2 Milestones	6
3.2.3 Sprints	6
3.2.4 Gantt chart	7
3.3 Effective Plan	7
3.3.1 Tasks	7
3.3.2 Milestones	7
3.3.3 Gantt chart	7
4 State of the art	11
4.1 Chatbots	11
4.1.1 History of Chatbots	11
4.1.2 Narrow	12
4.1.3 General	13
4.1.4 Artificial Intelligence Markup Language (AIML)	13
4.1.5 Information Retrieval (IR)	13
4.1.6 Deep Neural Networks (DNN)	13
4.1.7 Proactive	13

Contents

4.2	Word2Vec	13
4.2.1	What is Word2Vec	13
4.2.2	Gensim	13
4.2.3	Framworks	13
4.3	Word Embedding Alternatives	13
4.3.1	FastText	13
4.3.2	Glove	13
4.3.3	Word2Vec-f	13
4.3.4	Wang2vec	13
4.4	Sentence/Document Embedding Alternatives	13
4.4.1	Doc2vec	13
4.4.2	Skip-thought	13
4.4.3	Smooth Inverse Frequency	13
4.4.4	RNN	13
4.5	Datasets	13
5	Analysis	15
5.1	Natural Language Processing	16
5.2	Pipeline	16
5.3	Word2Vec	16
5.3.1	Bag of Words VS Skip-Gram	16
5.3.2	Dimensions	16
5.3.3	N-Grams	16
5.3.4	Epochs	16
5.3.5	Lemmatization	16
5.3.6	Normalization	16
5.3.7	Distance and Cosine Angle	16
5.3.8	Training	16
5.3.9	Retrain Model	16
5.3.10	Memory Issues	16
5.3.11	Analogies	16
5.3.12	Proverbs	16
5.3.13	Evaluation	16
5.3.14	Visual Representation	16
5.3.15	Benchmarks	16
5.3.16	CPU VS GPU	16
5.3.17	Datasets	16
5.4	Chatbot	16
5.5	Proactivity	16
6	Experiments & Results	17
6.1	Environments	17
6.1.1	Jupyter Notebook	17
6.1.2	Local Machine	17
6.1.3	Amazon Web Services	17
6.1.4	iColab-gpu2	17
6.1.5	CPU Dedicated Machine	17
6.2	Gensim Framework	17
6.3	Materials	17
6.4	Products	17

- 7 Discussion 19
 - 7.1 Next steps? 19
- 8 Conclusion 21
 - 8.1 Word Embedding: Word2Vec 21
 - 8.2 Framework: Gensim 21
 - 8.3 Word2Vec Chatbots 21
 - 8.4 Proactive Chatbots 21
- Appendix 25
 - .0.1 Appendix 25

Acknowledgments

If any

Acronyms

AGI

Artificial General Intelligence.

AI

Artificial Intelligence.

AIML

Artificial Intelligence Markup Language.

ANI

Artificial Narrow Intelligence.

ANN

Artificial Neural Networks.

DNN

Deep Neural Networks.

FAQ

Frequently Asked Questions.

IR

Information Retrieval.

ML

Machine Learning.

NLP

Natural Language Processing.

NLU

Natural Language Understanding.

Sci-Fi

Science Fiction.

Abstract

In the scope of this deepening project, and as the technology of NLP is in constant evolution, we will be focusing on the exploration of the word embedding algorithm Word2Vec, which is, at the beginning of 2019, commonly used as a foundation for DNN Chatbots. As a result to this project, the student is demonstrating what is the Word2Vec technology, its extensions, and its applications.

Keywords: Word Embedding, Word2Vec, Natural Language Processing (NLP), Natural Language Understanding (NLU), Machine Learning (ML), Data Engineering, Conversational Agent, Chatbot, Generic

Chapter 1

Introduction

Beginning of 2019, chatbots are everywhere but very limited to narrow tasks, and are, in most cases, sequences of if-else conditions resulting in a very weak Artificial Intelligence (AI). Indeed, hard-coded connections are requiring an infinite amount of human power to create generic Chatbots able to maintain a conversation at a human level. However, the progress in the field of ML is demonstrating that providing large corpora to an unsupervised algorithm is enough to maintain a passive conversation with users, which results into a shifting of the human power into data engineering. Multiple algorithms and technics are emerging monthly, which are demonstrating promising conversational performance improvement; however, they are all still narrow AI. Indeed, even if they are getting better at providing meaningful sentences, they are still not able to generalize all tasks linked to a conversation, such as, understanding the context, search and learn for missing information, initiate conversation in a meaningful manner, be intuitive, and more. The generalization of those features would allow a significant step forward into general Chatbots.

1.1 Aim of Study

In harmony with the author interest, the goal of this deepening semester project is to suggest and demonstrate strategic approaches as a premise to the AGI and getting a step closer to general Chatbots, which can initiate and maintain human-like conversations in a pro-active manner.

1.2 Scope and Study Borders

As a red line for this deepening project, the focus will be on the Word2Vec technology, from a research perspective. Indeed, this technology is seen as a foundation for the modern NLP and DNN Chatbots, which makes it an exciting vector of study about its current usage, its extensions, and potential evolution.

1.3 Industrial Interest

iCoSys, the Institut of Complex Systems at University of Applied Sciences and Arts at Fribourg, Switzerland, is interested into the result of this project as a study for their AI-News project, whose goal is to provide a chatbot as a tool to reader, to help them narrow their interests and deliver the right information. AI-News is in

Chapter 1. Introduction

collaboration with the Swiss Innovation Agency from the Swiss Confederation, and La Liberté, the daily newspaper from Fribourg.

Chapter 2

Questions

To help the student to find a red line to focus its research on, he was required work on the subject "What should be the initial questions to asks to start making Artificial General Intelligence (AGI) Chatbots" as preliminary study before the beginning of deepening project itself and to write down the outcome as a set of questions related to his interests and the field of AGI Chatbots.

2.1 Initial and Broad Questions

As a result to the preliminary study, the following question were produced. Please take into account that those questions were not meant to be answered as part of the project itself, but as part of the process of appropriation of the field of study.

- Is the Artificial Neural Networks (ANN) approach appropriate to represent the world?
- Can agents be made exclusively from a language?
- Can agents able to experience an environment?
- Is a narrative environment be enough to understand an environment?
- Is the language able to provide to an agent an understanding of the world?
- Is the knowledge of the language syntax enough to gain an understanding?
- Is the result of unsupervised learning enough to discover all nuances?
- Is the unsupervised learning sufficient to make sense to an environment?
- Is a descriptive explanation of the world be expressed in a language?
- Is the description good enough to catch all the nuances?
- Is the language good enough to explain?
- Can we augment or make a semantic language?
- Can we create a common symbolic language?
- Is the language multi-dimensional?
- How many dimensions are needed for a complex language?
- Is it possible to give a word equivalence to machines for human-specific words?
- Are all emotions describable into words?
- Are emotions altering language descriptions?
- Is an approximation of the real world enough to understand the environment?
- Would a the simulated world be a good approximation of the real world?

Chapter 2. Questions

2.2 Narrowed Questions

In a second time, the student was asked to narrow the initial questions above into potential fields of study.

- Common human-machine language
 - Is it possible to create a multi-dimensional human-machine language, which includes a common semantic, symbolic, and emotion definition.
 - Is it possible to create an abstract world for machines to understand human symbolic based on a real world, and define fundamentals for machine representation of the language.
- Machine intuition
 - Is it possible to provide to machines an human-like intuition (inside voice), which would help to keep a long term context and specialize in specific fields.
- Evaluate human-machine communication
 - Is it possible to provide a protocol to test the communication skills and machine understanding.

2.3 Potential Red lines

From the potential fields above, the following suggested red lines were proposed.

- How to verify and quantify a chatbot understanding?
- What is the premise to make chatbots general with today's technology?
- How chatbot can be proactive?
- How to simulate human-like intuition in chatbots?

2.4 The Deepening Project Question and Red line

Based on reflective work and discussions, the concluding red line and question for this deepening project is:

- What is Word Embedding and how is it useful for chatbots?

Chapter 3

Plan

3.1 Constraints

Timeframe: 15 weeks

Starting date: 18.02.2019

Ending date: 31.05.2019

3.2 Initial Plan

As an initial milestone for the deepening project, the student were required to create an initial plan, with the purpose to help the student and the teacher to visualise the project main red line.

3.2.1 Tasks

1. Initial research about general chatbots
2. Determine the project target
3. Play with the subject
4. Explore the Word2Vec methodology
5. Explore the Word2Vec extensions
6. Combine and test ANN algorithms with Word2Vec
7. Explore ANN algorithm topology for the chatbot
8. Analysis of the chatbot intuition with parallel algorithms
9. Analysis of a protocol to evaluate proactive chatbots
10. Profile-based initiatives
11. Analyze and experiment profile nurturing
12. Analyze and experiment with chatbot initiatives with no profiles
13. Overall improvements
14. Autonomous data gathering
15. Make suggestions
16. Determine possible continuation and future outcomes for the project

Chapter 3. Plan

3.2.2 Milestones

1. Initial deepening project plan and specification document
2. Basic multi-dimensional word embedding space
3. Basic conversational agent
4. Basic proactive chatbot
5. Deepening project report

3.2.3 Sprints

18.02.19 to 08.03.19 (3 weeks)

- Do the initial research about general chatbots
- Determine the project target
- Play with the subject
- **DELIVERABLE:** Plan and Initial Specification document

11.02.19 to 29.03.19 (3 weeks)

- Explore the Word2Vec methodology and its extensions
- Combine and test ANN algorithms with Word2Vec
- **MVP:** Basic multi-dimensional word embedding space

01.04.19 to 19.04.19 (3 weeks)

- Explore ANN algorithm topology for the chatbot
- Analysis of the chatbot intuition with parallel algorithms
- Analysis of a protocol to evaluate proactive chatbots
- **MVP:** Basic conversational agent

22.04.19 to 10.05.19 (3 weeks)

- Profile-based initiatives
- Analysis and experiment of the profile nurturing
- Analyze and experiment with chatbot initiatives with no profiles.
- **MVP:** Basic proactive chatbot

13.05.19 to 31.05.19 (3 weeks)

- Overall improvements
- Autonomous data gathering
- Make suggestions
- Determine possible continuation and future outcomes for the project
- **DELIVERABLE:** Report + Sources

3.2.4 Gantt chart

Figure 3.1 represents the visual gantt chart for the initial plan.

3.3 Effective Plan

As expected the initial plan served as an initial model, and evolved iteratively based on the student and teach feedback while exploring the subject.

3.3.1 Tasks

1. Initial research about general chatbots
2. Determine the project target
3. Set the initial plan
4. Make LaTeX report template
5. Explore the Word2Vec subject
6. Explore the Word2Vec algorithm
7. Build a Word2Vec model on the latest english wikipedia dump
8. Explore Word2Vec parameters
9. Explore Word2Vec analogies
10. Explore Word2Vec sentence generation
11. Explore visual representations of Word2Vec vectors
12. Explore Word2Vec applications with chatbots
13. Writing the report

3.3.2 Milestones

1. Initial deepening project plan and specification document
2. Basic Word2Vec Word Embedding Model
3. Conclusions Word2Vec based chatbots
4. Ideas to make chatbots proactive
5. Deliver the report

3.3.3 Gantt chart

Figure 3.2 represents the visual gantt chart for the effective plan.

Chapter 3. Plan

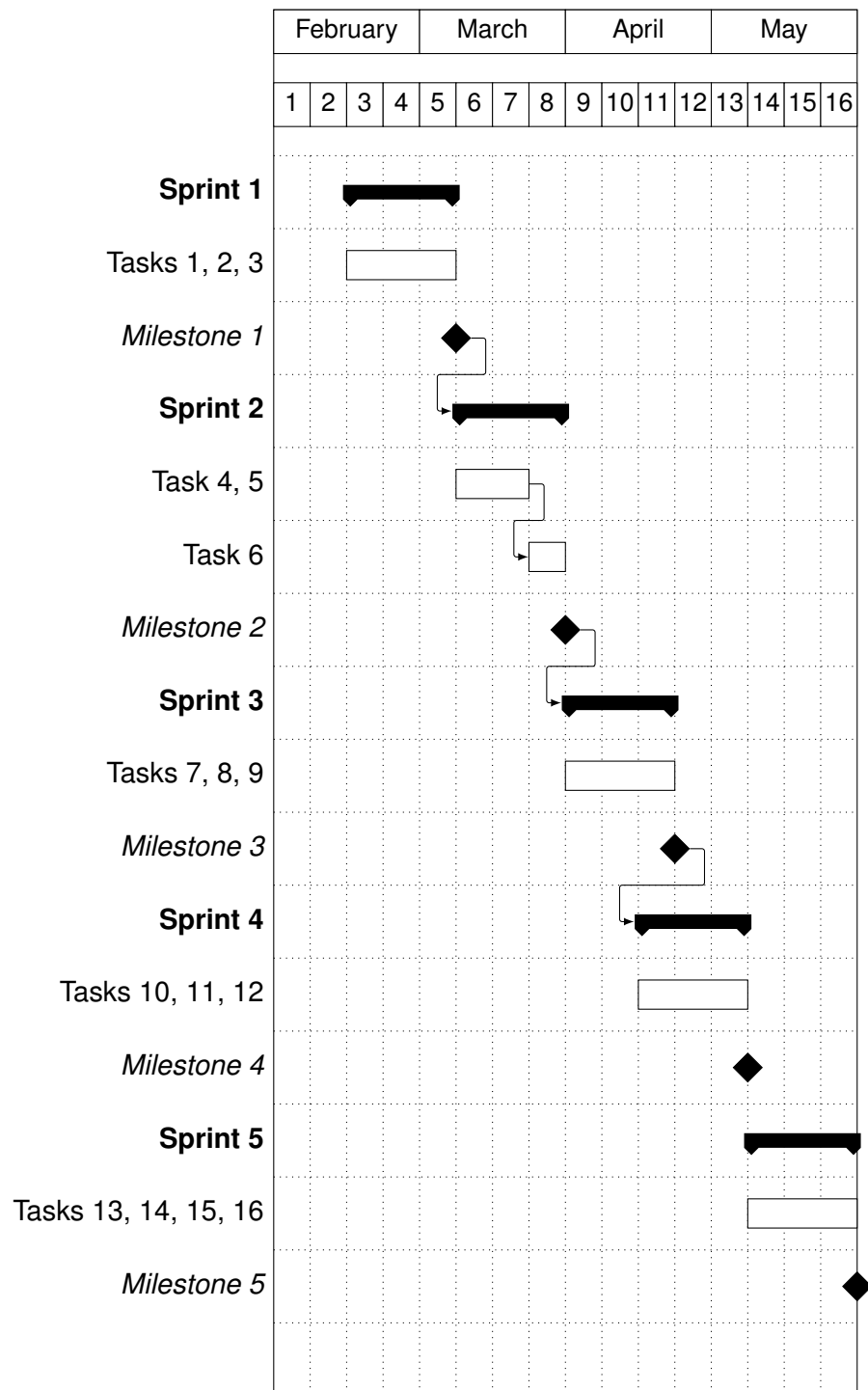


Figure 3.1: Initial Gantt Chart

3.3. Effective Plan

February					March				April				May			
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	

Figure 3.2: Effective Gantt Chart

Chapter 4

State of the art

4.1 Chatbots

From a user point of view, chatbots are trendy nowadays, big companies such as *Google* or *Apple* are pushing to make the technology mainstream. Even if not every lambda people understand the word Chatbot, they all have at least a mental representation of it. Indeed, they could call it Digital Assistant, Siri, Ok Google, etc, in the end they all understood the concept of an artificial intelligence narrowed to more or less human-like conversations.

4.1.1 History of Chatbots

When are they coming from? Not mentioning *Alan Turing* or *Joseph Weizenbaum*, who are considered as the fathers of AI and chatbots, would not be fair. Indeed, they forecasted in 1950, that computers will be able to use human-like communication and they proposed a test to distinguish humans from machines, called the Turing Test. Where a human is asked to talk to a masked entity, and determine if it is talking to a human or a computer. If the human cannot determine who is the computer, then the machine passed the Turing test, as seen on figure 4.1.

In 1966, Joseph Weizenbaum wrote Eliza, a computer program simulating a psychotherapist, which could be seen as one of the first well known attempt to make a chatbot passing Turing test. Note that due to technical restrictions, Eliza is not performing well at all time. As it is for today, it is possible to play with it at on a dedicated website. [3]

Since Eliza, a lot of progress has been made, indeed, to only cite a few noticeable chatbots: *Parry* (1972), *Jabberwacky* (1988), *Dr. Sbaitso* (1992), *A.L.I.C.E* (1995), *Smarterchild* (2001), *Watson* (2006), *Siri* (2010), *OK Google* (2012), *Alexa* (2015), *Cortana* (2015), Facebook Bots (2016), and *Tay* (2016), which where all part of the Chatbot history [2].

From IF-ELSE, AIML, up to ML with ANN and DNN, the improvement in the field of chatbots increased drastically over the years. At every iterations, the algorithms where becoming more sophisticated and better at using the human language, which are now called the field of the NLP and NLU.

Chapter 4. State of the art

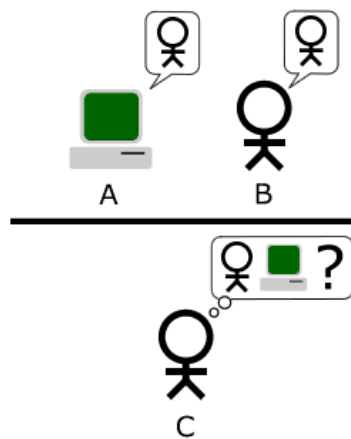


Figure 4.1: The "standard interpretation" of the Turing Test, in which player C, the interrogator, is tasked with trying to determine which player - A or B - is a computer and which is a human. The interrogator is limited to only using the responses to written questions in order to make the determination. [1]

4.1.2 Narrow

Once again, chatbots are almost everywhere nowadays. Indeed, it became a common tool for companies of any size to communicate with their customers and a toy for users. However, most of the time, Chatbots are not understood by their users, and is leading to a high level of frustration. Even if they are becoming increasingly mainstream and sophisticated, people doesn't realise their limits. Today's chatbots are often mistaken for AGI as it's seen in Science Fiction (Sci-Fi) and are expected to do much more than they are able to do. Indeed, we are making Artificial Narrow Intelligence (ANI) Chatbots, which implies a specialisation into a specific field.

We should not forget the main purpose of Chatbots, which to provide a conversational service to the user. However, its purpose can be applied in an almost unlimited amount of solutions. Health, Weather, Customer Service, Games, etc. And it could be declined in a text or vocal format.

Frequently Asked Questions (FAQ)

With the goal to answer specific questions, FAQ chatbots are probably the most common type of chatbots. Indeed, their communicational capacities are limited to pre-made sentence and a question answer database, which often results into in the best case scenario a perfect match, or in the worst case scenario something totally unexpected.

Sequential

Based on

Forwarder

Customer service

Learning

Examples

4.1.3 General

4.1.4 AIML

4.1.5 IR

4.1.6 DNN

4.1.7 Proactive

4.2 Word2Vec

4.2.1 What is Word2Vec

4.2.2 Gensim

4.2.3 Frameworks

4.3 Word Embedding Alternatives

4.3.1 FastText

4.3.2 Glove

4.3.3 Word2Vec-f

4.3.4 Wang2vec

4.4 Sentence/Document Embedding Alternatives

4.4.1 Doc2vec

4.4.2 Skip-thought

4.4.3 Smooth Inverse Frequency

4.4.4 RNN

4.5 Datasets

Chapter 5

Analysis

5.1 Natural Language Processing

5.2 Pipeline

5.3 Word2Vec

5.3.1 Bag of Words VS Skip-Gram

5.3.2 Dimensions

5.3.3 N-Grams

5.3.4 Epochs

5.3.5 Lemmatization

5.3.6 Normalization

5.3.7 Distance and Cosine Angle

5.3.8 Training

5.3.9 Retrain Model

5.3.10 Memory Issues

5.3.11 Analogies

5.3.12 Proverbs

5.3.13 Evaluation

5.3.14 Visual Representation

5.3.15 Benchmarks

5.3.16 CPU VS GPU

5.3.17 Datasets

5.4 Chatbot

5.5 Proactivity

Chapter 6

Experiments & Results

6.1 Environments

6.1.1 Jupyter Notebook

6.1.2 Local Machine

6.1.3 Amazon Web Services

6.1.4 iColab-gpu2

6.1.5 CPU Dedicated Machine

6.2 Gensim Framework

Errors Memory allocation with multi-core. The problem is occurring during the merge of the cores. Indeed, my current machine has 128GB ram, and the dataset weights about 16GB in the memory, and each core during merging is processing at least the same amount, plus the processed informations.

```
2019-03-25 08:31:18,867 : INFO : PROGRESS: pass 0, dispatched chunk #34 = documents up to #70000/4614519, outstanding queue
size 31
Exception in thread Thread-1:
Traceback (most recent call last):
  File "/usr/lib/python3.5/threading.py", line 914, in _bootstrap_inner
    self.run()
  File "/usr/lib/python3.5/threading.py", line 862, in run
    self.target(*self.args, **self.kwargs)
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 366, in _handle_workers
    pool._maintain_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 240, in _maintain_pool
    self._repopulate_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 233, in _repopulate_pool
    w.start()
  File "/usr/lib/python3.5/multiprocessing/process.py", line 105, in start
    self._popen = self._Popen(self)
  File "/usr/lib/python3.5/multiprocessing/context.py", line 267, in _Popen
    return Popen(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 20, in __init__
    self._launch(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 67, in _launch
    self.pid = os.fork()
OSError: [Errno 12] Cannot allocate memory
```

Figure 6.1: Error 1

6.3 Materials

6.4 Products

Chapter 6. Experiments & Results

```
2019-03-25 08:31:18,867 : INFO : PROGRESS: pass 0, dispatched chunk #34 = documents up to #70000/4614519, outstanding queue
size 31
Exception in thread Thread-1:
Traceback (most recent call last):
  File "/usr/lib/python3.5/threading.py", line 914, in _bootstrap_inner
    self.run()
  File "/usr/lib/python3.5/threading.py", line 862, in run
    self.target(*self.args, **self.kwargs)
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 366, in _handle_workers
    pool.maintain_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 240, in _maintain_pool
    self.repopulate_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 233, in _repopulate_pool
    w.start()
  File "/usr/lib/python3.5/multiprocessing/process.py", line 105, in start
    self._popen = self._Popen(self)
  File "/usr/lib/python3.5/multiprocessing/context.py", line 267, in _Popen
    return Popen(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 20, in __init__
    self._launch(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 67, in _launch
    self.pid = os.fork()
OSError: [Errno 12] Cannot allocate memory
```

Figure 6.2: Error 2

Chapter 7

Discussion

7.1 Next steps?

Chapter 8

Conclusion

- 8.1 Word Embedding: Word2Vec**
- 8.2 Framework: Gensim**
- 8.3 Word2Vec Chatbots**
- 8.4 Proactive Chatbots**

Lausanne, May 27, 2019

Romain Claret

Bibliography

- [1] Bilby. *File:Turing Test version 3.png*. [Online; accessed 26-May-2019]. 2008. URL: https://commons.wikimedia.org/wiki/File:Turing_Test_version_3.png.
- [2] Futurism, LLC. *The History of Chatbots Infographic*. [Online; accessed 26-May-2019]. 2016. URL: <https://futurism.com/images/the-history-of-chatbots-infographic>.
- [3] Michal Wallace & George Dunlop. *Eliza, the Rogerian Therapist*. [Online; accessed 26-May-2019]. 1999. URL: <http://psych.fullerton.edu/mbirnbaum/psych101/Eliza.htm>.

Appendix

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

.0.1 Appendix

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

