

POLITECNICO DI MILANO
Master of Science in Computer Science and Engineering
Dipartimento di Elettronica, Informazione e Bioingegneria



[Title of the Thesis]

Supervisor: [Name]
Co-supervisor: [Name]

M.Sc. Thesis by:
[Name], matriculation number [nnn]
[Name], matriculation number [mmm]

Academic Year 200N-200N+1

About this template

With this template I want to give you some input on how to structure your thesis if you develop your thesis with me in Politecnico di Milano. Next to the pure structure, which you should reuse and adapt to your own needs, the document also contains instructions on how to approach the different sections, the writing and, sometimes, even the work on your thesis project itself. Sometimes you will also find boxes like this one. These are meant to provide you with explanations and insights or hints that go beyond the mere structure of a thesis.

I hope this template will help you do the best thesis ever, if not in the World, at least in your life.

Florian Daniel
October 12, 2017

Disclaimer: Sometimes I may make statements that are general, if not over-generalized, personal considerations, or give hints on how to do work or research. Be aware that these are just my own opinions and by no way represent official statements by Politecnico di Milano or its community of professors. If something goes wrong with your thesis or presentation, you cannot refer to these statements as a defense. You are the final responsible of what goes into your thesis and what not.

Acknowledgements: The original template for this document was not created by me. I would love to acknowledge the real creator, but I actually do not know who it is. The template has been passed on to me by a former student, who also didn't know the exact origin of it. It was circulating among students. However, to the best of my knowledge at the time of writing, it seems that Marco D. Santambrogio and Matto Matteucci may have contributed at some point with considerations on structure and funny citations. Both were helpful and enjoyable when preparing this version of the template. I will be glad to add more precise acknowledgements if properly informed about the origins of this template.

Supervisors and co-supervisors

If the supervisor is internal to Politecnico di Milano (a professor or researcher), then on the first page use "Supervisor" plus the titles "Prof." and "Dr." for professors and researches, respectively. If the work was co-supervised by someone else, refer to him/her as the "Co-supervisor." If the work was supervised by someone external to Politecnico di Milano, use "External supervisor" for the external supervisor plus "Internal supervisor" for the internal supervisor that mandatorily must co-supervise the work with the external supervisor.

Optionally, here goes the dedication.

Abstract

The abstract is a small summary of the thesis. It tells the reader in few words (up to one/one and a half page of total text) everything he/she needs to understand:

- ☐ the *context* of the work (e.g., chatbots),
- ☐ the specific *problem* approached by the thesis (e.g., the development of personal bots by non-programmers),
- ☐ if applicable, clearly state the *research questions* you would like to answer (e.g., “is it possible to enable non-programmers to do X using A?”),
- ☐ the three/four *core aspects of the proposed solution* (e.g., use pre-defined rules, use machine learning, assisted development, etc.),
- ☐ the *concrete outputs* produced by the thesis (e.g., a state of the art analysis, a conceptual/mathematical model, an application, middleware or API, an empirical study with/without users, etc.), and
- ☐ the *findings and conclusions* that one can draw from the evaluation of the approach (e.g., that under some very specific conditions non-programmers are indeed able to implement own chatbots effectively using the proposed technique).

Checklists

Now and there I propose checklists with items, such as the one just above this box. They are meant for you to check if you included all the content that is relevant and that should be included, in order to make your text complete. When reading your thesis, I will look for all these items.

Writing style

This is a M.Sc. thesis. It's neither Facebook nor Twitter nor an email. This is going to be an official document with legal value that will decide on the final mark of your yearlong university career and perhaps even on your future work perspectives. So, you surely don't want to be judged badly because of grammar errors, flawed/wrong vocabulary or superficial layout and/or text structure. It is a must that what you write is always *correct* content- and language-wise (no false statements or claims, no language mistakes), *readable* (no sentences that cannot be understood) and targeted at the *average-skilled reader* (professors, but also your own colleagues).

Plagiarism

This is a M.Sc. thesis. It's neither Facebook nor Twitter nor an email. This is going to be an official document with legal value that will decide on the final mark of your yearlong university career and perhaps even on your future work perspectives – yes, I plagiarized myself here a little bit. So, you surely don't want to copy/paste material from scientific articles, online resources, books, and similar without adequately acknowledging the holders of the respective intellectual property rights. If you do so, it is a must that you properly *cite* each source where you take text or inspiration from. It is fine to do so – actually, citing someone is a compliment! – but it becomes a crime if the source is not cited. Not only M.Sc. titles but also Ph.D. titles have been withdrawn for fraudulent “reuse” of others' intellectual property. Be aware that Politecnico di Milano, like most higher educational institutions that issue university degrees or scientific publishers, may use specialized software to automatically detect plagiarism.

Sommario

Here goes the translation into Italian of the abstract. If the thesis is written in Italian, no translation into English is needed. Hence, one of the following must be checked:

- ☐ Thesis written in *English*, properly proofread translation needed
- ☐ Thesis written in *Italian*, no translation needed, chapter omitted

Acknowledgements

If you would like to thank somebody for given support, this is the right place to do so.

Contents

Abstract	I
Sommario	III
Acknowledgements	V
1 Introduction	1
1.1 Context: [topic]	1
1.2 Scenario and Problem Statement	3
1.3 Methodology	4
1.4 Contributions	5
1.5 Structure of Thesis	7
2 State of the Art	9
2.1 [Topic one]	9
2.2 [Topic two]	10
2.3 Summary	10
3 Data Collection	11
3.1 Tools	11
3.1.1 Tweepy	11
3.1.2 Botometer	12
3.1.3 Hoaxy	12
3.2 Datasets	12
3.2.1 Caverlee-2011	12
3.2.2 Cresci-2017	12
3.2.3 Varol-2017	13
3.2.4 BotBlock	14
3.3 Varol clustering	14
3.4 Collection	18
3.4.1 NSFW	19

3.4.2	News-spreader	19
3.4.3	Spam-bots	21
3.4.4	Fake-followers	21
3.4.5	Genuine	22
3.5	Data visualization	23
4	Features engineering	29
4.1	Baseline	30
4.2	Missing values filling	32
4.3	Descriptive features	33
4.4	Intrinsic features	34
4.5	Extrinsic features	36
4.6	Image features	38
4.7	final feature vector	39
4.7.1	User features	39
4.7.2	Tweets features	39
5	Model selection	41
5.1	Baselines	41
5.1.1	Random Forest	42
5.1.2	Logistic Regression	42
5.1.3	K-Nearest Neighbors	43
5.1.4	Support Vector Machine	44
5.1.5	Comparison and baseline selection	45
5.2	Multiclass classifier	48
5.2.1	Dataset	49
5.2.2	Model	49
5.3	Binary Classifier	53
5.3.1	Dataset	53
5.3.2	Model	54
5.4	Text classifier	57
5.5	Stacking meta-classifier	57
5.5.1	Genetic algorithm	60
5.5.2	Logistic Regression	65
5.6	Prediction pipeline	70
6	Web application - BotBuster	73

7	Implementation and Evaluation	75
7.1	Implementation	76
7.2	Evaluation	76
7.2.1	Design of Evaluation	76
7.2.2	Metrics	77
7.2.3	Results	78
7.2.4	Discussion	78
8	Conclusion and Future Work	79
8.1	Summary and Lessons Learned	79
8.2	Outputs and Contributions	79
8.3	Limitations	80
8.4	Future Work	81
	References	83
A	User Manual	85
B	Dataset	87

Chapter 1

Introduction

The introduction is one of the core chapters of your thesis. It expands what has already been said in the abstract with additional details on the content and contribution and on the structure of the thesis. It is meant to introduce the reader to the work he/she will be reading in the rest of the document and, most importantly, to get the reader curious about reading on, knowing more about your work.

1.1 Context: [topic]

This thesis is about describing the work you are doing in your final thesis project. You have been working on it for months, and nobody knows the work better than you do. This is great and exactly how things should be: by doing your thesis project you became an expert – if not *the* expert – in this specific field of research and/or technology.

But attention: being the expert is also dangerous when it comes to explaining others what you did and why you think you did a great work that deserves attention (I give it for granted that your work does so). There are only very few people around you (your supervisor and possible co-supervisor, some friends, maybe someone else) who are as expert as you are in this topic. So, if you start in a full-impact fashion to tell that you implemented an extraordinarily cool, new algorithm to solve X, or that you discovered this extremely surprising finding Y, or that you mathematically proofed that Z, etc. (you got it), your reader will not understand anything. Therefore, before talking about what you actually did, you need to introduce the reader to the context of your work, provide the necessary core definitions that are needed to understand the terminology you will be using in the rest of the thesis (if it's not standard IT terminology).

Therefore:

- ☐ Tell the *research area(s)* your work/project focuses on. If you are doing your thesis with me, likely candidates of research areas are Web Engineering, Data Science, Crowdsourcing, Service-Oriented Computing, Business Process Management.
- ☐ Tell possible *sub-areas* that are more specifically related to what you are doing. Again, if you are doing your thesis with me, likely candidates of sub-areas are chatbots, social knowledge extraction, business process matching/modeling, quality control in crowdsourcing, etc.
- ☐ Make the *heading* of your context section self-explaining by substituting “[topic]” in heading 1.1 with the sub-area most relevant to your work. It should read like “Context: quality control in crowdsourcing” or similar.
- ☐ If needed, introduce some *key definitions* (no need to introduce everything here, but be sure that the introduction does not use terminology the reader may not be familiar with). For instance, if you are working on chatbots, this is definitely a term that needs to be introduced here; it’s not yet commonly known but it’s crucial for the understanding of the rest of the thesis and introduction.
- ☐ Use *examples* to make definitions and ideas concrete and clear.
- ☐ Throughout, make *references* to the relevant literature.

Use of tenses and pronouns

Writing a thesis is writing a scientific document like scientific articles or research publications. There are two conventions that are usually applied in this kind of publications (admittedly, they may seem somewhat odd if not used to):

First, the most used tense is the *simple present*. The thesis is meant to describe a piece of work, from problem statement, to the conception of a solution, its implementation and evaluation. Yet, it’s not a novel about your life, and it’s not meant to provide a chronological story about what you did and didn’t do. Content is presented in an order that is most effective to convey its message, not in time order. In this spirit, it’s much more effective to say “in order to get result A, first we do X, then we do Y and then Z,” instead of saying “in order to get result A, we did Y after having done X, then we went on doing Z.” The order of actions, their interconnections, inputs and outputs already tell the dependency – if properly described. Most of the times, the most effective way to describe a solution or

methodology only becomes clear after trial and error. It's enough to explain the result, not how you got there chronologically.

Second, the *pronoun* used to talk about the own work is "our" (work). That is, it is custom to say "we" instead of "I," even if you are writing your thesis alone. However, don't forget about all the people that helped you get there: your supervisor, co-supervisor, colleagues, etc. This may sound strange at the beginning, but, at the other hand, using "I" too often risks to convey the impression that you are self-focused and egoistic, which is never good.

1.2 Scenario and Problem Statement

Now that the reader got the general context of your work and has an intuition of the problem you will be solving in the rest of the thesis, it's time to be clear about which *specific problems* your thesis project is going to solve. One way of doing so is by describing a *scenario* (a description of a real situation, with all its actors, roles, tasks, instruments, etc.) that provides evidence that there are one or more real problems right now that, with the current technology and understanding of the domain, are hard to solve or not solvable at all. If instead the problem(s) can be solved already, it should be evident from the scenario that this is possible only at a prohibiting cost or with unsatisfying guarantees on the quality of the result or not within useful time for the target user.

It's important that the scenario is written in such a way that the reader, after reading it, agrees with you that the problem you are focusing on is a relevant one, one that deserves being studied and solved. Consider that if you convince the reader here that your thesis is needed (after all, that's what this section is about), he/she will be very open to possible solutions and happy to see how you solve it. If instead you fail to convince the reader – let me be harsh – the whole rest of your thesis is useless in the eyes of the reader. This is the worst outcome you want.

Conclude this section by explicitly stating which of the problems evident in the scenario you are approaching. Don't raise false expectations! Never ever tell the reader there are five core problems and then solve only two of them in the thesis, without telling upfront that this is what you intended to do in the first place. As soon as you list problems, the reader wants to see a solution, unless you stop him/her immediately from thinking so by telling that out of the described problems you focus on a subset only, usually because this subset is already a huge research and development problem in its own.

In summary:

- ☐ Describe a *real scenario* that provides evidence of *real problems*.
- ☐ Convince the *reader* that the problems need to be solved.
- ☐ Use an *illustration* or *figure* to help the reader understand.
- ☐ If possible, provide *references* to literature that backs your assessment of the problem.
- ☐ Provide a clear *problem statement* that summarizes what came out of the scenario and your specific focus.

1.3 Methodology

Fixed the problem(s) you want to approach, you can approach it/them in thousands of different ways. Your way is just one of the thousands, and the reader may have (and very likely will have) a very different intuition of how to solve the problem(s) you just pointed out. So, clarify how you intend to proceed:

- ☐ Tell if you follow an existing *methodology* or not; if yes, name it and provide a reference to literature, if available. For example, Design Science [?] is a likely methodology to cite here.
- ☐ Tell which of the following *procedures, techniques, methods* you use in your work and for which purpose (put them also into the right order, so that their application or use makes immediate sense to the reader):
 - ☐ *Systematic literature review, survey*
 - ☐ *Statistical hypothesis formulation and testing*
 - ☐ *Software prototyping*
 - ☐ *Iterative development*
 - ☐ *Participatory design*
 - ☐ *Performance evaluation*
 - ☐ *Comparative studies*
 - ☐ *User studies*
 - ☐ *Expert interviews*
 - ☐ *Simulation/emulation*

- ☐ *Live experiments*
 - ☐ *Case studies*
 - ☐ *Mathematical theorem proving*
 - ☐ *Mathematical modeling*
 - ☐ *Pseudocode*
 - ☐ *Graphical modeling* (e.g., UML, ER)
 - ☐ *Model-driven development*
 - ☐ *Automatic code generation*
 - ☐ ...
- ☐ Tell if you use some special *software instruments* that help you in your work. We are of course not talking about Word or Google Search. Perhaps you can tell that you used R for data analysis or some specific modeling instrument for automated code generation or simulation.

1.4 Contributions

Now that the reader knows what you want to solve and how you intend to proceed, you can anticipate the contributions your thesis makes to the state of the art. Attention, a thesis project may produce lots of different *outputs* (e.g., a software prototype, a set of registrations and transcripts of interviews, datasets collected during experiments) and *contributions* (e.g., a demonstration that some software solutions solves a given problem under well defined conditions, a formal proof that some property holds, empirical evidence that something works as expected). The former are all the artifacts produced throughout the work. The latter refer to *new knowledge* (if you are doing a full thesis) or the most important, *final output* (if you are doing a tesina). Sometimes, outputs and contributions overlap, but not necessarily.

Typical contributions are (multiple choices may apply to your thesis):

- ☐ A *systematic literature review* of the state of the art providing evidence for some argument
- ☐ The design of a *model* (mathematical, graphical, algebraic, etc.) describing how to solve a real world problem in a reusable fashion
- ☐ The drawing of *conclusions* (findings) from the analysis of a dataset describing some physical or virtual phenomenon

- ☐ The implementation of a *software prototype* solving a real world application problem
- ☐ The design of a *language* (textual, graphical) enabling others to solve own problems or to solve them easier
- ☐ *Formal proofs* of correctness, completeness or other properties of the proposed models or theorems
- ☐ *Objective evidence* from empirical studies (e.g., performance analyses or simulations) that demonstrate that the proposed prototype or solution works / works better than existing software or solutions that solve the same/similar problem(s)
- ☐ *Subjective evidence* from user studies or expert interviews backing the claims of viability of the proposed problem or solution/artifact
- ☐ A reasoned *argumentation*, e.g., based on a detailed case study, supporting the viability of the proposed problem or solution/artifact

Thesis vs. Tesina

Let me spend some words on the difference between these two. Before that, however, it is important to clarify the very purpose of your final project, be it a thesis or a tesina (a small thesis). The purpose of it is giving you the possibility to show that, after years of attending classes and giving exams, you are also able to *apply* the knowledge you acquired during your studies. In short, it's all about you showing that you are *mature*. Mature from a knowledge perspective, mature from an application perspective, mature from a work/teamwork perspective, mature from an ethical perspective.

It is common that a thesis project is not very well defined in its beginning and that even the supervisor does not really know how to approach a given problem or which problem to focus on in the first place. This may even be annoying to you, but attention: there is no intention behind it. Your supervisor is not withholding information from you to test you or to see if you get something. It's just the nature of real *problem solving*. If things were clear from the beginning, there wouldn't be any problem! Fledging out the problem and agreeing on a solution and methodology is a core part of you demonstrating your maturity – if not the most important one. *How* you proceed from the inception of the thesis idea to the final solution is as important as *what* you find and/or produce in the end.

This being said, a *thesis* in Politecnico di Milano usually requires you to make a contribution to the literature (the so-called state of the art). Making a contribution – from a science point of view – means creating new *knowledge*, that is, finding something that nobody knew before, demonstrating a property that nobody showed

before, improving the performance of a given system with a new algorithm, and similar. For a thesis, it is therefore not enough to produce a perfectly engineered solution. It is key that you also demonstrate, provide empirical evidence or proof that your solutions performs as claimed. Well, for a *tesina* this last demonstration is usually not required, and the focus is on the engineering of the solution. In addition, perhaps in the case of the *tesina* the solution to be engineered is also less complex then for a thesis, but this depends on the context and on how you want to measure complexity.

1.5 Structure of Thesis

Here you explain the structure of the thesis, so that the reader knows how to read it. Consider that not every reader wants to read through the whole thesis to find some specific information. Actually, only few will do so (your supervisor and co-supervisor, and the possible reviewer for sure). Many more will just leaf through it and look for specific types of information (e.g., the context of your work, your findings, how you implemented something, which technologies you used). It is your duty to accommodate them all. How? By telling them how your thesis is structured.

Therefore, in this section you provide a brief description (2-3 sentences) for *each* chapter that follows this introduction. Use an itemized or numbered list to structure the text, like this:

- ☐ Chapter 2 introduces the state of the art and...
- ☐ Chapter 3 provides...
- ☐ ...

Structuring text

Besides telling the reader how the content of your thesis is organized into chapters, it is important that you master some basic text structuring techniques. To organize your text there are lots of instruments you can use: chapters, sections, sub-sections, paragraphs, itemized lists, numbered lists, code examples, figures, images, screen shots, captions below figures, tables, and so on. Use them all! Don't write text without structure. Never.

Be aware that the structure of your text, that is, how you present your work, conveys a lot of information about how well you actually understand what you are writing about, how much you care about being clear and helping your reader understand, and how much value you give yourself to your own thesis. A well

structured presentation of content that the reader can understand and agree with is a huge plus in this respect. Text that lacks proper paragraphs, does not use lists where needed, etc. is a minus and also much harder to read (think about how much a well structured text can help you go back ten pages and find concepts you know you read about compared to a text that comes without an easy to memorize formatting and structure). When writing, think about some of your textbooks. Since you are doing an engineering degree, I'm sure these are textbooks that make exemplary use of the different formatting instruments available.

Chapter 2

State of the Art

This chapter discusses the state of the art that is relevant for your own work. What does that mean? It means that it provides the reader with all the relevant references he/she may need to know in order to understand better three things: (i) the context of your work, (ii) the problem and the need for a solution, and (iii) the value of your contribution. You achieve this by citing works or scientific papers that solved the same or similar problems in the past. Citing does not just mean adding a references to the bibliography and printing a number here; it means you tell the reader about the merits and possible demerits of each of the references you feel relevant. Of course, doing so requires you to first read each reference and, most importantly, to understand it. There should be lots of references in this chapter.

It is advisable that you structure the chapter into sections in function of the topics you treat. If you do so, before starting with the first section of the chapter, explain the reader how you structure your discussion in one paragraph.

- ☐ *Read* relevant literature and or *test* related software or tools.
- ☐ *Summarize* your reading.
- ☐ Provide correct *references* (the bibliography in the end of this document).

2.1 [Topic one]

...

2.2 [Topic two]

...

2.3 Summary

Close the state of the art chapter with some words that connect the discussion of the references to your thesis. Pay attention that the reader understands why you discussed the works/topics you discussed and how they are related to what you do.

- ☐ Show that in the state of the art the *problem* you want to solve has not yet been solved or not been solved in an as efficient / effective / easy to use / cost-saving fashion as you target with your work.
- ☐ If your work has similarities with some *specific references*, point them out here and explain why these are particularly important to you. Perhaps you started your investigation from the outputs of a specific paper or you want to improve the performance of an algorithm studied earlier; it's good to mention this here.
- ☐ Attention: this is not yet the place where to anticipate *your solution*. You may give hints, but it's too early to make a comparison between your work and the state of the art, as the reader does not yet know anything about your work. This discussion can go into the final chapter.

Chapter 3

Data Collection

In this chapter we will present all the available datasets containing bot accounts, with all the tools and methodologies used to collect new data. The final dataset contains:

- ⇒ Data from existing datasets
- ⇒ Data collected with different approaches
- ⇒ Hand-labeled data

3.1 Tools

Different tools were used in order to both collect the data and to enrich existing ones. This stage was essential to gather additional features. Here we present all the instruments involved in this section.

3.1.1 Tweepy

Tweepy is a python wrapper for the Twitter API. Two main methods were used to collect all the default features of users and tweets:

Method	Input	Output
API.get_user	[id/user_id/screen_name]	User object
API.user_timeline	[user_id][, count]	list of Status object

A User object contains all the features that describe the user's profile and his utilization of Twitter. A Status object contains the details of a single tweet, such as the full text, number of retweets and replies.

3.1.2 Botometer

Botometer [3] checks the activity of a Twitter account and gives it a score based on how likely the account is to be a bot.

Using their dataset they tested most classification algorithm and highlighted Random Forest as their final classifier, since it gained the highest accuracy. Its accuracy was 98.42% and 0.984 F1 measure[4].

Botometer provides API to check Twitter accounts genuinity.

3.1.3 Hoaxy

Hoaxy is a tool that visualizes the spread of articles online. Articles can be found on Twitter, or in a corpus of claims and related fact checking.

Hoaxy only checks the news sources and compares them with a list of unreliable URLs.

3.2 Datasets

3.2.1 Caverlee-2011

This dataset is composed by content polluters, detected by a social honeypot, and legitimate users sampled from Twitter. For each content polluter, they save their 200 most recent tweets, their following and follower graph, and the temporal and historical profile informations.

In order to collect genuine users, they randomly sampled about 20.000 Twitter ids and monitored them for three months, to see if they were still active and not suspended by the social platform moderation service [4].

3.2.2 Cresci-2017

Cresci Dataset is composed by genuine accounts and bots. In this dataset there is a deeper differentiation for the bots, which are precisely labeled according to different categories.

Dataset	Description	#Users	#Tweets	Year
genuine accounts	verified accounts that are human-operated	3,474	8,377,522	2011
social spambots #1	retweeters of an Italian political candidate	991	1,610,176	2012
social spambots #2	spammers of paid apps for mobile devices	3,457	428,542	2014
social spambots #3	spammers of products on sale at Amazon	464	1,418,626	2011
traditional spambots #1	training set of spammers used by Yang[2]	1,000	145,094	2009
traditional spambots #2	spammers of scam URLs	100	74,957	2014
traditional spambots #3	automated accounts spamming job offers	433	5,794,931	2013
traditional spambots #4	automated accounts spamming job offers	1,128	133,311	2009
fake followers	accounts inflating followers of other accounts	3,351	196,027	2012

- **Genuine accounts** are those users who correctly answered to a simple question, posed in natural language, so they represent accounts with no automatization.
- During Rome majoral election in 2014, one of the candidates used a set of automated accounts to publicize his policies. These accounts were gathered to be part of the **Social spambots #1**.
- **Social spambots #2** are accounts that promotes mobile app, using popular hashtags for months.
- **Social spambots #3** promotes products on sale on Amazon, by tweeting products URL and descriptions.

All these accounts were manually checked to verify their automated nature.

- The **Traditional spambots #1** dataset is the training set used in [2].
- **traditional spambots #2** are users that mention other users in tweets containing scam URLs. They usually invite users to claim a prize.
- **Traditional spambots #3** and **traditional spambots #4** are bots that continously tweet job offers.
- **Fake followers** are account involved in increasing popularity of other users. In order to collect them, they bought followers from fastfollowerz.com, intertwitter.com and twittertechnology.com. [9]

3.2.3 Varol-2017

This dataset contains a list of Twitter accounts, labeled as bots (1) or humans (0).

The construction of the Varol dataset starts with the identification of a representative sample of users, by monitoring a Twitter stream for 3 months, starting in October 2015. Thanks to this approach it is possible to collect data without bias; in fact other methods like snowball or breadth-first need an initial users set. During the observation window about 14 million user accounts has been gathered. All the collected users must have at least 200 tweets in total and 90 tweets during the three month observation (about one tweet per day).

Using the classifier trained on the honeypot dataset in [4], they computed the classification scores for each of the active accounts. Then the samples

were grouped by their score and 300 accounts from each bot-score decile were randomly selected. The 3000 extracted accounts were manually labeled by some volunteers. They analyzed users profile, friends, tweets, retweets and interactions with other users. Then they assigned a label to each user. Of course the final decision is conditioned to personal opinion[6].

3.2.4 BotBlock

Botblock (<https://github.com/dansarie/Botblock>) is a Twitter block list containing the user ids of a large number of known porn bot accounts. They are mainly used to aggressively market porn sites.

3.3 Varol clustering

At first try, we wanted to understand if different kinds of bots were easy to distinguish, using their profile features only. er, we didn't know what kinds of bots populate Twitter for the most. So, an unsupervised approach could have helped us to highlight different categories. We relied expectations on clustering techniques, hoping to get a solid help in automatizing the labeling process of the data.

We used the Varol dataset[6], that contains a plain list of bots and humans. It was not possible to use all the data, because we needed to scrape from the web browser all the possible features. Some of the listed accounts were already been deleted, so, for this work, we had to consider only those accounts that were still active.

The first step was the knee-elbow analysis based on a hierarchical clustering with single linkage and euclidean distance. The data must have been preprocessed and cleaned.

Here we illustrate what kind of preprocess operations were performed for that purpose:

feature	type	preprocess operation
id	int	delete
name	str	replace with len(name)
screen_name	str	replace with len(screen_name)
statuses_count	int	—
followers_count	int	—
friends_count	int	—
favourites_count	int	—
listed_count	int	—
url	str	replace with hasUrl (0/1)
lang	str	one hot encoding
time_zone	str	one hot encoding
location	str	one hot encoding
default_profile	int	replace with hasDefaultProfile (0/1)
default_profile_image	boolean	boolean to int (0/1)
geo_enabled	boolean	boolean to int (0/1)
profile_image_url	str	delete
profile_use_background_image	boolean	boolean to int (0/1)
profile_background_image_url_https	str	delete
profile_text_color	str	delete
profile_image_url_https	str	delete
profile_sidebar_border_color	str	delete
profile_background_tile	boolean	boolean to int (0/1)
profile_sidebar_fill_color	str	delete
profile_background_image_url	str	delete
profile_background_color	str	delete
profile_link_color	str	delete
utc_offset	int	delete
is_translator	boolean	boolean to int (0/1)
follow_request_sent	int	delete
protected	boolean	boolean to int (0/1)
verified	boolean	boolean to int (0/1)
notifications	boolean	delete
description	str	replace with hasDescription (0/1)
contributors_enabled	boolean	boolean to int (0/1)
following	boolean	delete
created_at	str	string to int (year)



Figure 3.1: Hierarchical Clustering Dendrogram (truncated to the last 20 merged clusters)

Since we didn't know how many categories of bots were listed in this dataset, the first step consisted in understanding which was the optimal number of clusters to look for. To achieve that goal, we applied *hierarchical clustering*. In figure (3.1) you can see the dendrogram of the algorithm.

In order to select the optimal number of clusters, we plotted the knee-elbow figure (3.2). It shows the variation of WSS (within cluster sum of squares) and BSS (between cluster sum of squares) as the number of clusters increase.



Figure 3.2: Hierarchical Clustering - knee-elbow

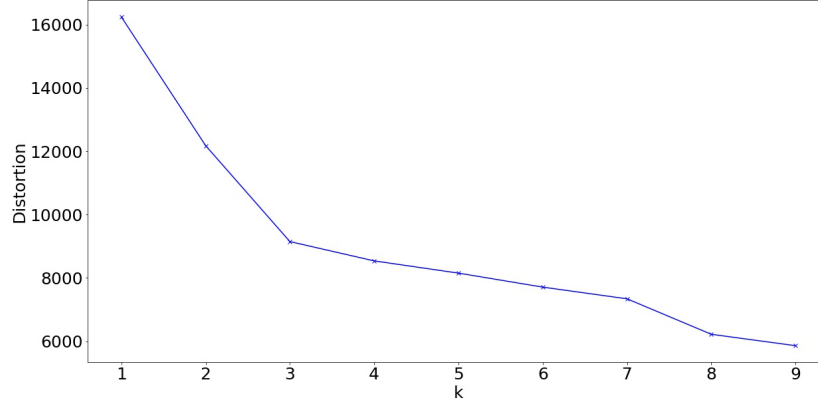


Figure 3.3: The Elbow Method showing the optimal k

It is clear that there are no well-defined elbows or knees, both curves seems to be "smooth", so it were hard to pick a reasonable k (number of clusters) for the algorithms to come.

Then we tried another approach. We applied the *K-means* algorithm and we plotted the elbow method (3.3). In this figure there is an elbow between $k=3$ and $k=4$, so the most accurate solution were represented by four clusters.

As this process came to an end, we manually inspected the resulting clusters.

cluster	size
cluster 1	82
cluster 2	648
cluster 3	2
cluster 4	16

Cluster 2 contains most of the samples, while the others have fewer elements. We observed the Twitter profile of all the elements belonging to cluster 1, 3, 4 and a small sample of profiles for cluster 2.

Unfortunately there was no correlation among accounts in the same cluster, so this technique didn't seem to fit the speeding up of the labeling process, nor to create useful features for a classifier.

3.4 Collection

The clustering approach did not help us, but the manual inspection of the clusters allowed us to get in touch with some existing bots, making us understand which categories of bots are most common on the social network. In particular we detected 4 main classes:

- ⇒ NSFW bots
- ⇒ News-spreader
- ⇒ Spambots
- ⇒ Fake-followers

We started with a hand-labeling of the Varol dataset [6]. For each account we analyzed its profile and tweets and we assigned it a label according to the categories we identified. We have faced some unexpected behaviors among bot accounts, that didn't fit the above-mentioned categories. In these cases, they were temporarily signed as "general purpose".

We also found genuine users, who we thought that had been incorrectly added to the dataset. This task culminated with the collection of the following bots:

category	labeled account
NSFW	31
news-spreader	71
spambots	418
fake-followers	5
general purpose	63
genuine	104

"*General purpose*" accounts are sometimes bots with no goal, they aim to emulate human behavior and often they were recognizable just because their description informs other users about their own nature.

Sixtythree users were not enough to represent a class and it was not possible to find a large list of those accounts who act like them, so we added all these ids to the "genuine" group. Even if this choice brought some noise to our data, that allowed us to provide our data more heterogeneity.

"*NSFW*" accounts are only thirtyone elements, anyway the problem of pornography is a known issue on Twitter. In [6] they clusterize users too.

”These bot clusters exhibit some prominent properties: cluster C0, for example, consists of legit-looking accounts that are promoting themselves (recruiters, porn actresses, etc.)” [6]. This kind of users are often banned by Twitter, so it is likely that the accounts that we were not able to scrape, used to belong to this category. Therefore we firmly believed that obtaining further accounts of this class was fundamental.

Finally it is clear that even fake-followers were few, since they were not considered in the Varol research, but they are important in [9], so we decided to expand this category too.

All these samples were not enough to train a classifier, hence we needed to collect more data. We perform this task by focusing on one category at a time.

3.4.1 NSFW

Not safe for work is a tag used on internet to mark all that URLs, web pages, e-mails that contain nudity, intense sexuality, profanity or violence. In particular we wanted to collect a specific sub-category: the pornbots. In order to collect them, we used the BotBlock dataset.

BotBlock contains thousands of pornbot ids. We wanted to gather about 6000 samples, an amount that would have been enough for the final dataset, with regards to its balance. Since they were sorted according to their creation date, we shuffled the whole list. Then, using the Twitter API, we looked for accounts that weren’t deleted yet. We needed to scrape profile features and tweets, so we couldn’t consider banned accounts. The user list was initially shuffled to allow us to collect users with different ages. This emerged as a very useful setp, because we gathered both more long-lived accounts and more extreme accounts (which probably have shorter lives). We finally obtained 6903 users and 198378 tweets.

3.4.2 News-spreader

Many bots on Twitter are news-spreader. The goal of these users is to spread politics, sports or actuality news. Often their behavior is not harmful, they just retweet statuses from newspapers accounts. However, there are users created to diffuse fake news. In the last few years Twitter has been used to boost politics propaganda. During elections or political campaigns, ad hoc accounts are created to divulgate specific political idea.

As a recent study highlighted, about the 80% of these ”pre-elections bots” are still alive [8]. We think that part of our news-spreader dataset includes some of those accounts.



Figure 3.4: Collection of news-spreading bots, approach 1

We started gathering these ids by exploring Hoaxy. We used two different approaches. The first way (in figure 3.4) consisted in collecting the twenty top popular fake news, for each month, in the last two years. We performed this task using the Hoaxy APIs. Thanks to this service, we obtained all the tweet ids that have spread the considered claim. With the official Twitter APIs we collected all the users involved in this spreading activity. We finally passed all these accounts to the Botometer API, since many of the retrieved users were humans.

We set a threshold, in order to classify a user as a bot. That threshold is 2.3 due to the willing of including some false positives in our data, increasing the heterogeneity in their behaviors and the challenge level for the classifiers. We think that a high intra-homogeneity among classes could lead the models to perform well on the training data, but worse over unseen ones.

The second approach just consisted in collecting the most popular news-spreaders according to Hoaxy. In consistency with the mentioned threshold, profiles with a Botometer score lower than 2.3 were still discarded.

Finally we checked every profile added to this dataset and removed all

that users who didn't tweet enough statuses to be included in this class. This hand-made analysis made us know that there are no bots who only spread fake-news. Usually they tweet a lot of verified news and some fake ones, to keep their credibility. We reached 3590 accounts and 333699 tweets.

3.4.3 Spam-bots

As seen before, spambots were already collected in [9]. Authors allowed us to access to their dataset, so we obtained the spambots list by sampling their data. Due to homogeneity reasons, we needed to perform scraping again, since we needed different features compared to the available ones. We selected users from:

- ⇒ traditional spambots 1
- ⇒ social spambots 2
- ⇒ social spambots 3

We chose this categories because they contain the most popular kinds of spambots, that are the ones who advertise products, services or mobile applications. We ignored "*social spambots 1*", since they are italian news-spreaders and "*traditional spambots 3 and 4*", since we retrieved enough job-offer spammers during the hand-made labeling. If we had stored too many bots of this category, we would not have been able to generalize on generic spambots. Finally we gathered 4943 accounts and 458809 tweets.

3.4.4 Fake-followers

The collection of this class was quite easy. We initially performed scraping of the data collected by the Cresci research [9]. Many of these accounts had already been banned, so we could not collect their features. In order to enrich our dataset, we created a new Twitter account (figure 3.5). Then we bought fake followers from two different services:

- ⇒ instakipci.com/
- ⇒ rantic.com/buy-legit-twitter-followers/

Instakipci provides low-quality followers. Usually they have no tweets, no followers and a they have a lot of followings.

Rantic, on the other hand, ensures more realistic followers. They seem to have a real network of friends and, sometimes, they tweet too.

By using both services, we gathered a more miscellaneous dataset. We collected 6307 users and 41683 tweets.



Figure 3.5: Collection of fake-followers bots

3.4.5 Genuine

Finally we needed genuine accounts. We used again the Cresci dataset [9] and we filled it with all the Varol users labeled as humans. We again performed scraping on the existing accounts and we collected 3661 users and 263240 tweets in total.

3.5 Data visualization

As the collection of the data was completed, we explored our final dataset. A first look (3.6) shows us how many user accounts we collected (y axis) for each class and how many tweets we could scrape (diameter). It is easy to observe that fake-followers and nsfw bots have less tweets than the others, while news-spreaders have a lot of tweets, but we collected less profiles.

category	target id
NSFW	0
news-spreader	1
spambots	2
fake-followers	3
general purpose	4



Figure 3.6: users amount and tweets

Then we plotted the heatmap of the correlation matrix (3.7). We wanted to understand if some feature was more useful to predict the correct target. This plot suggested us that there were no feature highly correlated to the target.

Moreover in figure (3.8) it is possible to see the distribution of the missing values. This is fundamental for the features engineering step.



Figure 3.7: heatmap

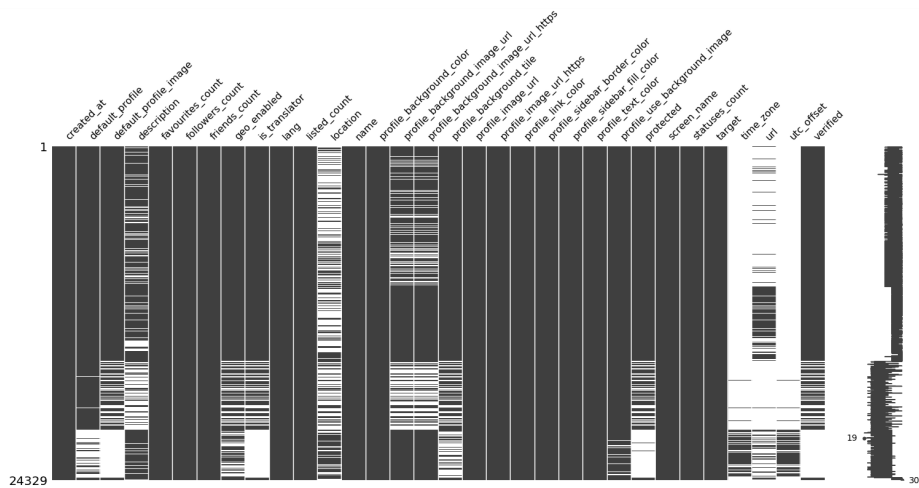


Figure 3.8: heatmap



Figure 3.9: boxplot statuses_count

Finally we performed three further analysys. With the heatmap we could not detect which features were more important. Anyway,during the hand-labeling step, we understood that some of these features were instead very usefull to identify a few classes of our dataset. These attributes are *"followers_count"*, *"friends_count"* and *"statuses_count"*. For each of them we plotted a boxplot. It is a method used to represent groups of numerical data through their quartiles. In Figure 3.9 we can analyze the statuses count for each target. We collected up to 100 tweets for each user, so this chart is limited to 100. Here we can notice an interesting behavior: *"news-spreaders"* and *"spambots"* are the classes with more tweets, while *"fake-followers"* have less statuses. By reflecting on the gols of the bots, this result is exactly what we expected to see. In fact *"fake-followers"* don't need to tweet, they just need to exist, while other types of bots have to publish many statuses, to draw attention on their contents.



Figure 3.10: boxplot friends_count

In Figure 3.10 and 3.11 we performed the same analysis considering the "friends_count" and "followers_count" features. We limited the charts to 4000 and 2000 to keep them understandable. These two figures show us that "news-spreaders" usually have a bigger network, while "fake-followers" just follow few users. "News-spreaders" network may depend on the popularity of the news media. Other categories are more balanced and their differences can be attributed to the data and not to a different behavior between them. For the genuine accounts we need to make an extra observation. Of course there are a lot of human users on Twitter and they have very different behaviours. In our dataset we included profiles with an average popularity (3.11), so they are neither VIPs nor fresh users without followers. Furthermore they have a smaller network of friends (3.10), which is in the average of the other categories. Finally they actively use twitter (3.9) since we were able to collect about 100 tweets for each of them.



Figure 3.11: boxplot followers_count

Chapter 4

Features engineering

This chapter can be seen as one of the most important of the whole project.

We wouldn't have hit such performances, if it wasn't for feature engineering. We had a large pool of models to pick for our purpose, and we tried different assets for them, but the difference were made with the reasoning behind the construction of the final feature vector.

Twitter APIs provides us two kinds of features: the user attributes and the tweet attributes. We knew that user attributes weren't enough to infer on targets, so we started planning how to include tweet informations and enhance our data with them. This was the bulk of the work, but it helped us to catch characteristic behaviors of some user. We created many features. Some of them are descriptive, like the lenght of strings (name, description ecc) or the count (miminimum, maximum, average) of other aspects like hashtags per tweet and tweet's lenght. Features describing the tweeting activity (frequency, how often a tweet contains a media or a url or it is just a retweet) have been considered too. Other kind of features are more behaviour-oriented, like the monotony of different tweets of the same user, while others are oriented to the most used words in tweets. Finally we also added features related to image analysis.

The choice over the amount of tweets that would have been considered was a trade-off between performance and prediction speed.

We finally chose to retriive up to the latest 100 tweets for each user, because further material led us to a slower, but equivalent, prediction over test samples.

At the end of this stage, the resulting - and final - feature vector will include 38 features.

4.1 Baseline

In this section we analyse the complete set of default profile features and which kind of pre-processing operation we applied. With this default set we trained several classifiers to define some baselines for the upcoming work and this allowed us to evaluate the improvements made by the features engineering step.

Some features are ready to be used in a classifier, while other ones need to be pre-processed, in order to allow them to be more expressive. We wanted to rely on user features only for this experiment, in order to have a large improvement margin to exploit, once we would have gone deeper in the study.

During the data exploration stage, we identified the most and least meaningful attributes to trace a baseline, so we started from that. We simply tried to improve and to homologate the features highlighted in the previous chapters.

We faced a lot of missing values as well as non-numeric ones. Even if the goal was to build a raw model with semi-raw data, we needed feasible and manageable attributes to work with.

Her we list all the pre-processing operations applied to each feature belonging to the ones provided by the method *get_user()*, of the official Twitter APIs.

feature	type	preprocess operation
id	int	delete - useless feature
name	str	delete - non-numeric feature
screen_name	str	delete - non-numeric feature
statuses_count	int	—
followers_count	int	—
friends_count	int	—
favourites_count	int	—
listed_count	int	—
url	str	replace with hasUrl (0/1)
lang	str	delete - non-numeric feature
time_zone	str	delete - too many missing values
location	str	delete - too many missing values
default_profile	int	delete - too many missing values
default_profile_image	boolean	boolean to int (0/1)
geo_enabled	boolean	delete - too many missing values
profile_image_url	str	delete - non-numeric feature
profile_use_background_image	boolean	boolean to int (0/1)
profile_background_image_url_https	str	delete - non-numeric feature
profile_text_color	str	delete - non-numeric feature
profile_image_url_https	str	delete - non-numeric feature
profile_sidebar_border_color	str	delete - non-numeric feature
profile_background_tile	boolean	boolean to int (0/1)
profile_sidebar_fill_color	str	delete - non-numeric feature
profile_background_image_url	str	delete - non-numeric feature
profile_background_color	str	delete - non-numeric feature
profile_link_color	str	delete - non-numeric feature
utc_offset	int	delete - too many missing values
is_translator	boolean	delete - too many missing values
follow_request_sent	int	delete - relative feature
protected	boolean	delete - too many missing values
verified	boolean	delete - too many missing values
notifications	boolean	delete - relative feature
description	str	replace with hasDescription (0/1)
contributors_enabled	boolean	delete - too many missing values
following	boolean	delete - relative feature
created_at	str	delete - useless feature

Features processed as "delete - relative feature" are those ones related to the user who performed the scraping. So we didn't need them.

4.2 Missing values filling

Features with few missing values was not deleted from the dataset, but we needed to fill that fields. In this section we analyze how we performed this task for each features.

- ⇒ **default_profile_image:** Thanks to the figure 3.8 we could see that all the missing values was at the bottom, in particular, all of them was in tuples with target 3 or 4. In order to uderstand the behaviour of this feature, we printed its values count for each indicted target.

target	3
value	mean
0	2868
1	228

target	4
value	mean
0	181
1	19

In both cases the value "0" is more frequent then "1", so we filled al the missing values with the mode (0).

- ⇒ **profile_background_tile:** As for "default_profile_image", all the missing values belonged to target 3 and 4. We used the same approach and we obtained:

target	3
value	mean
0	3086
1	99

target	4
value	mean
0	1347
1	147

In this case we decided to fill these fields with the mode (0). Since most of the data are 0, this choice allowed us not to dirty the dataset.

⇒ **profile_use_background_image**: All the missing values are still in the last two classes.

target	3
value	mean
0	12
1	4983

target	4
value	mean
0	25
1	3246

This data is really unbalanced, so filling the null fields with the mode (1) is still the better solution.

4.3 Descriptive features

In order to enrich our attributes and to provide support to our algorithms, we decided to add some descriptive "meta" features, such as synthesis statistics and counters.

Their purpose is to describe the tweets in a statistical way, adding ranges and means to the attributes provided by the official APIs.

Each of these new values were been added to the users feature vector, in order to append new informations for each account. Here is the list of this first 18 brand new features, introduced by our work:

feature	description
avg_len	average lenght of the tweets (words)
max_len	lenght of the longest tweet (words)
min_len	lenght of the shortest tweet (words)
avg_ret	average amount of retweets (by other users) per tweet
max_ret	highest amount of retweets (by other users) on a tweet
min_ret	lowest amount of retweets (by other users) on a tweet
avg_fav	average amount of favourites (by other users) per tweet
max_fav	highest amount of favourites (by other users) on a tweet
min_fav	lowest amount of favourites (by other users) on a tweet
avg_hash	average amount of hashtags involved in tweets
max_hash	highest amount of hashtags involved on a tweet
min_hash	highest amount of hashtags involved on a tweet
freq	amount of tweets per day (up to 100)
ret_perc	percentage of retweets, made by the user, over its retrieved tweets
media_perc	percentage of media content incorpored in tweets
url_perc	percentage of URL links placed inside tweets
quote_perc	percentage of quotes, made by the user, over its retrieved tweets

4.4 Intrinsic features

Due the multiclass nature of our dataset, it was impossible to rely on the descriptive meta features only.

We faced the need of better capturing some behaviours, that could have helped us distinguish between targets.

We spent a lot of time analysing Twitter timelines by ourselves. This was one of the most useful phases of our work.

Indeed, we have learnt a lot about bots acting like humans on the social platform. One thing that was easy to notice was the monotony, in terms of words or URLs involved in tweets, met with Spam-bots, as well as the opposites, for Genuine accounts or Fake-Followers.

We tried to encapsule this distinctive behaviour by adding two intrinsic features to the training vector, along with the descriptive ones.

How to portray such monotony?

We though about different approaches, like complex sentiment analysis or entity recognition, but then, we chose to rely on a weighting technique and the euclidean distance.

We looked inside every retrieved tweets for each user, then we encoded each

of them with TF-IDF weighting.

Every term (word) of every tweet was represented by a numeric weight, according to TF-IDF.

This weighting formula is a combination of Term Frequency (TF) and Inverse Document Frequency (IDF).

$$TF_{i,j} = \frac{n_{i,j}}{|d_j|}$$

The term frequency factor counts the number n of the i_{th} term inside the j_{th} document (the tweet, in our case), dividing it by the length of the latest, in order to give same importance to both short and long collections of texts.

$$IDF_i = \log \frac{|D|}{|\{d : i \in d\}| + 1}$$

Where d is the document (tweet).

The inverse document frequency factor aims to highlight the overall magnitude of the i_{th} term in the collection which it belongs. The collection D , in our work, is represented by all the gathered tweets of the examined user.

$$(TF - IDF)_i = TF_{i,j} \times IDF_i$$

After the encoding process, we wanted to map the resulting vectors into an euclidean space, in order to compute the distance of each weighted text, from the total centroid of the collection.

We decided to add each user a measure of the average intra-distance of its tweets.

In order to achieve this goal, we relied on the WSS metric used in K-means clustering, but trying to soften its magnitude. We didn't want huge ranges in our features, in order to minimize the normalizations along the process.

The resulting formula for this brand new attribute is the following:

$$IntraDistance(U) = \frac{1}{N} \sum_{\mathbf{x} \in U} \|\mathbf{x} - \boldsymbol{\mu}\|^2$$

Where N is the number of tweet for user U , x is the encoded tweet and μ is the centroid of the tweet collection for that user.

This formula, as well as the whole process, has been used both for tweet words and URLs inside of them. The resulting 2 features are

feature	type
tweet_intradistance	float
url_intradistance	float

We expected low values for Spam-bot accounts in both features, as well as we expected the opposite for Genuine and Fake-Follower accounts.

4.5 Extrinsic features

Once we have modelled the personal twitting actions, in terms of words and links dissimilarity, we needed to look for those parameters that could be compared with all the users in our dataset. We wanted the users to get out of their shells, and to match their timelines with each others.

Once again, sentiment analysis came to our mind. We found lots of paid or limited services that could have only partially supported us during this stage.

We couldn't think about implementing our own semantic analysis, as it is meant to be, due to the effort and time it would had taken.

For simplicity sake, we had the idea to look for the most meaningful words in tweets, that are common to all the users belonging to the same class.

We tried this approach, hoping to find a robust help in separating topics among targets.

The idea was to build five partially-non-overlapping dictionaries, one for each class of users, containing the most popular words used by them in their retrieved tweets, stripped of stopwords. Each dictionary is ordered according to the occurrence of each word, in order to have a proper ranking for the terms and to give each of them a score.

We gathered the 300 most common words for each category of accounts, in order to average about 250 terms, once the overlapping elements would have been discarded. Considering more than 300 words led us to irrelevant performance boost, instead it required lot of time to process a single user, so we decided to stop early.

We have listed 5 dictionaries:

dictionary	size
NSFW_main_words.csv	209
spam_bots_main_words.csv	230
news_spreders_main_words.csv	237
fake_followers_main_words.csv	233
genuine_main_words.csv	223

The overall scores are normalized, so that the most common word, for each class, is associated with a unitary weight, the least common one with a value similar to zero (it depends on the final amount of words included in the dictionary).

In terms of representation, this *keywords score* is built with 5 different values, one for each class.

The dictionaries aren't totally disjointed with each others. We decided to strip, from each list, only the top 50 words belonging to the other collections, since the most relevant scores reside in that zone of the ranking. We didn't want to give the user high scores, for the same word, but in different categories.

If some user hit a word that is placed inside more than a dictionary, we wanted to let the relative weights speak and assign scores to each of the targets that contain that word, knowing they are less relevant, as they don't fall in the highest positions.

So we have

feature	type
NSFW_words_score	float
spam_bots_words_score	float
news_spreders_words_score	float
fake_followers_words_score	float
genuine_words_score	float

When we process a new user and infer on him, we scan every words of every tweets and match them with all the dictionary. Every time that we hit a listed word, we assign the user the score of that word.

For instance, If the word we are comparing matches with the most frequent for NSFW accounts (the top position in its dictionary), we update the *NSFW_words_score* of our user, summing 1 to its current value for that feature.

We expected to capture patterns about the choice of the words involved in tweets, for each of our categories. This extrinsic attributes revealed themselves as very useful, lately.

4.6 Image features

One of the main issues of our work was to make the NSFW class different, in terms of classifiers vision, with respect of Spam-bots.

According to our user-based, descriptive and intrinsic features, both classes act in a similar way: they spam similar links, have high tweet frequency and their contents are often repeated.

What could be the difference?

We started with "blind" classifier, that was the main problem. Our feature vector lacked in visual components. A Spam-bot could have been detected as a NSFW, as its tweets involved media and URLs. We needed to go deep and actually see what kind of media were broadcast.

This is the reason for presence of the upcoming features.

We've found a small versatile project on GitHub for NSFW detection. The projects involves a pre-trained TensorFlow neural network for image recognition, that looks for adult and violent content inside pictures. It assigns the media a probability to be not safe for work.

We decided to scan our dataset with this new component and to give a NSFW score to both profile pictures and tweets.

For time complexity reasons, we couldn't imagine to scan all the retrieved tweets for each user and to look for embedded media. We limited the process to the latest ten tweets.

Obviously, the prediction time has been affected by this new preprocessing stage, but we think that this number of pictures analysed (up to eleven), makes the generalization duration still reasonable.

The project creators claim to reach 98.2% of accuracy in their final test, so we thought that the bias introduced with this new features would have been under control and that it wouldn't compromise the final results of our models.

The brand new attributes that helped us in better sperating NSFW class are:

feature	type
NSFW_profile	float
NSFW_avg	float

4.7 final feature vector

At the end of this process we obtained a feature vector composed of 38 elements, 12 based on the user and 26 based on his tweets.

4.7.1 User features

Features
age
default_profile
description_len
favourites_count
friends_count
followers_count
listed_count
profile_use_background_image
name_len
screen_name_len
statuses_count
url

4.7.2 Tweets features

Descriptive features
freq
min_fav
avg_fav
max_fav
min_hash
avg_hash
max_hash
min_len
avg_len
max_len
min_ret
avg_ret
max_ret
media_perc
quote_perc
ret_perc
url_perc

Intrinsic features

tweets_intradistance
url_intradistance

Extrinsic features

NSFW_words_score
news_spreaders_words_score
spam_bots_words_score
fake_followers_words_score
genuine_words_score

Images features

NSFW_profile
NSFW_avg

Chapter 5

Model selection

In this chapter we will show the choices and stages behind the final model. Starting from baseline models, we enhanced the chosen classifiers with hand-crafted features coming from the last chapter.

We saw and studied the performance improvements with validation approaches, and this phase led us to our current solution.

The result involves three models:

- ⇒ a first Random Forest classifier that has been used to provide an early filter on the separation between Genuine accounts and Bots
- ⇒ a second Random Forest that gives a classification among the five studied categories
- ⇒ a Naive Bayes classifier, used over the same classes of the second Random Forest, but it reads and labels the users, according on their tweets only

All the above mentioned algorithms were combined with a stacking ensemble methods, after considering different possibilities.

5.1 Baselines

The choices explained in this section were made at the same time of the ones listed in the Baseline section of the last chapter.

This is, basically, the same stage of the above-mentioned, but in a model-driven perspective. The features involved are the ones described in that chapter, but we started from that base, to try different classifiers over it. Each classifier has been fitted with the entire dataset, but considering only baselines features.

Furthermore, no parameters tuning has been applied, in order to minimize the results of our baselines classifier, with their standard settings.

5.1.1 Random Forest

Random forest is an ensemble learning method used in classification tasks and prediction ones as well.

The algorithm builds several *decision trees* and the resulting output is provided by the mode of the predictions coming from the estimators in the forest.

Each decision tree is trained on a subset of the original data, formed by sampling with replacements the whole training set. They share the same splitting criterion, in order to build subtrees, which is the entropy: Every tree computes the Information Gain of each feature, which is the difference, in terms of entropy, between the information gained on the data D , before splitting on the attribute X , and the one gained after the split, which provides n subsets of D .

$$InformationGain(X) = Information(D) - Information_X(D)$$

where

$$Information(D) = -p_1 \log p_1 - \dots - p_n \log p_n$$

and

$$Information_X(D) = \frac{|D_1|}{|D|} Information(D_1) + \dots + \frac{|D_n|}{|D|} Information(D_n)$$

The attribute providing the highest InformationGain, against the others at the same level of the tree, is chosen to perform a split.

The feature set considered by each tree is a random subset of the original pool.

Due to its ability to face overfitting and to the feature importance ranking that it can provide, this tool is often preferred over other models belonging to the same category.

The advantage of preventing overfitting usually comes with a slower prediction time, because it needs enough estimators for this task. But, for our purpose, there were enough estimators to face the variance problem without affecting the generalization speed.

5.1.2 Logistic Regression

Logistic regression is a common statistical model, that uses a sigmoid function to map the output of a linear regression on a normalized score, giving

the probability, for each sample, to belong to the positive class, given its features and a weighting vector:

$$P(\hat{y}_i = +1 | \vec{x}_i, \vec{w}) = \frac{1}{1 + e^{-\vec{w}h(\vec{x}_i)}}$$

Where \hat{y}_i is the predicted target, over the i_{th} sample, \vec{x}_i is the feature vector of that sample, \vec{w} represents the weighting vector that has to be learned and h is the activation function of the linear regression.

Logistic Regression searches for the weighting vector that matches the highest likelihood and, in order to do that, it minimizes a cross-entropy error function, provided by the negative log of the likelihood:

$$\mathbf{L}(\vec{w}) = -\ln \prod_{i=1}^n P(\hat{y}_i = +1 | \vec{x}_i, \vec{w})$$

In multiclass tasks, there are two possible approaches to face the problem:

- ⇒ a more general *softmax* function to replace the logistic sigmoid, which assigns the probability, for the i_{th} sample, to belong to the class C :

$$P(\mathbf{C}_i | \vec{x}_i, \vec{w}) = \frac{e^{-\vec{w}h(\vec{x}_i)}}{\sum_{j=1}^n e^{-\vec{w}h(\vec{x}_j)}}$$

- ⇒ "One-vs-Rest" method, which for each class, it builds a model that predicts the target class against all the others.

We decided to stick with the default settings of the libraries involved, so OvR was the approach used for the baseline.

5.1.3 K-Nearest Neighbors

K-Nearest Neighbors is an instance-based model used for classification, regression and pattern recognition. It is considered as a lazy learning algorithm, because all the computation is deferred until the prediction phase. When it performs a classification over a new point, it looks for the K nearest samples in the training set, according to a chosen metric, and it assigns, to the unseen sample, the mode of the targets of the retrieved neighbors.

The choices to make are the ones regarding the number K of neighbors to consider, the weights to assign to them and the metric to calculate the distance with. We used the default settings for the metric (*Euclidean distance*) and for the weighting technique (*uniform*), but we chose to consider

10 neighbors, because the automatic setting was $K = 5$, which is the number of our possible targets. We chose a K that is large enough to make the model not too sensible to outliers, and restricted enough to sharpen the classes boundaries.

We first normalized the training data and then we fitted the algorithm on them, in order to simplify the distance computations.

5.1.4 Support Vector Machine

Support Vector Machine is a smart way to do instance-based learning. It can be seen as a generalization of the weighted KNN algorithm, with an arbitrary and feasible *kernel function*, instead of the more generic dot product.

It can be summarised with a support vector $\tilde{\mathbf{x}}$ (a subset of the training set), a weighting vector $\tilde{\mathbf{w}}$ for them and a **kernel** $K(x, x')$ (a similarity function).

In order to make it work properly, three choices must be made:

- ⇒ a proper kernel, which is often selected according to experience and domain knowledge of the problem. We wanted to make things simple in this stage, so we used the default kernel function, which is the Radial Basis Function:

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

with σ as a free parameter

- ⇒ the weights \vec{w} , which are obtained by maximizing the margin that splits the records belonging to different classes. Each samples are mapped into a space, thanks to what is known as the *kernel trick*. The "trick" helps a linear classifier to work on a non-linear problem, applying the kernel function in the prediction phase.

This process highlights the boundary that separates the points belonging to different classes. SVM aims to draw the boundary for the classes, in order to maximize the "margin" formed between the closest points that have different targets

- ⇒ the support vector \vec{x} , which comes as a consequence of choosing weights

Since we were still facing a multitarget problem, the binary nature of SVM must had been adapted to our needs. We decided, once again, to stick with the default setting for non-binary classifications, in order to have only raw baselines to compare.

The multitarget classification is handled with "One-vs-One" approach. It considers all possible pairwise binary classifiers and so it leads to $\frac{N(N-1)}{2}$

individual binary classifiers, where N is the number of the classes in the problem.

In comparison with "One-vs-Rest" approach, "One-vs-One" is less sensitive to an imbalanced dataset, but it's more computationally expensive than the other, which only builds N binary classifiers. Despite our choices over methods and parameters weren't accurate in this stage as they were in the other ones, we decided to stick with this setting for SVM, because otherwise it would have led us to an irrelevant algorithm, in comparison with the above-mentioned.

5.1.5 Comparison and baseline selection

The selected baseline models were tested with a holdout approach at first, then with a crossvalidation method. We built a Confusion Matrix for each model, in order to bring out goodness indices for each class, such as *True Positive* (TP), *False Positive* (FP) and *False Negative* (FN). The evaluation metrics considered are *Precision*, *Recall* and *F1 score* and they work on the mentioned indices.

$$\Rightarrow Precision = \frac{TP}{TP+FP}$$

It measures the proportion of positive identifications, for a given target, that was actually correct

$$\Rightarrow Recall = \frac{TP}{TP+FN}$$

It measures the proportion of actual positive classifications that was identified correctly

$$\Rightarrow F1score = \frac{2(Precision \times Recall)}{Precision + Recall}$$

It calculates the harmonic mean of the previous metrics

Every metric is adapted to fit a multiclass problem. For each class, it has been computed this set of measures, and then they were averaged without weights (macro average), in order to not take label imbalance into account.

The results are pretty similar between the two methods, however we used the one coming from crossvalidation to select the model to build.

Holdout evaluation

The holdout stage is performed separating the samples in the dataset into training and test subsets. The splitting process is randomized, and we decided to use the 75% of the data for the training set and the 25% for the test set. This choice is a little bit different from the most common one, which builds the training set with 2/3 of the whole data, because we didn't dispose

of a huge amount of records, so we preferred this ratio and then trying an other validation method for comparison. Here we list the algorithm and their parameters, as they were written according to the Scikit-learn library for Python, their confusion matrix and their scores:

⇒ *RandomForestClassifier*(*n_estimators* = 10, *criterion* = 'entropy')
Confusion matrix:

		Predicted class				
		NSFW	NS	SB	FF	GEN
Actual class	NSFW	1654	23	14	4	68
	NS	30	716	31	3	81
	SB	6	27	1223	3	53
	FF	15	4	6	1212	3
	GEN	67	74	45	3	718

Precision: 0.895

Recall: 0.894

F1 score: 0.894

⇒ *LogisticRegression*(*fit_intercept*=True, *max_iter*=100, *penalty*='l2')
Confusion matrix:

		Predicted class				
		NSFW	NS	SB	FF	GEN
Actual class	NSFW	1310	176	78	40	159
	NS	26	676	54	1	104
	SB	25	60	947	221	59
	FF	167	27	11	1032	3
	GEN	118	295	172	8	314

Precision: 0.675

Recall: 0.685

F1 score: 0.673

⇒ *KNeighborsClassifier*(*n_neighbors*=10)
Confusion matrix:

		Predicted class				
		NSFW	NS	SB	FF	GEN
Actual class	NSFW	1512	39	50	82	80
	NS	46	674	28	14	99
	SB	45	24	1077	38	128
	FF	88	4	94	1018	36
	GEN	138	69	132	86	482

Precision: 0.769

Recall: 0.762

F1 score: 0.765

⇒ `SVC(kernel='rbf', decision_function_shape='ovo')`

Confusion matrix:

		Predicted class				
		NSFW	NS	SB	FF	GEN
Actual class	NSFW	1761	0	1	0	1
	NS	861	0	0	0	0
	SB	1068	0	244	0	0
	FF	449	0	0	791	0
	GEN	907	0	0	0	0

Precision: 0.468

Recall: 0.364

F1 score: 0.321

Crossvalidation

This approach is based on repeated holdouts. It is performed by splitting the whole data in K non-overlapping folds, leading to K different holdout evaluations. The results for each step are stored and the final evaluation is given by the mean of the K evaluations. For each evaluation, one fold is used for testing, the other ones for training the models. A common practice is to set $K = 10$ and thus averaging 10 different evaluations. This method is also known as *K-fold crossvalidation*. We used a stratified approach, which takes care about keeping the labels balanced on each fold.

Due the need of performing ten steps, it is computationally more expensive than a simple holdout validation. In our case, it was feasible, in terms of speed, because of the models complexity and the data amount.

The obtained scores are also more meaningful, with regards to holdout, because they are less sensitive to "lucky" or "unlucky" splits.

Here is the results for every baseline model:

⇒ *RandomForestClassifier*(*n_estimators* = 10, *criterion* = 'entropy')

Mean precision: 0.895

Mean recall: 0.891

Mean f1 score: 0.890

⇒ *LogisticRegression*(*fit_intercept*=True, *max_iter*=100, *penalty*='l2')

Mean precision: 0.729

Mean recall: 0.704

Mean f1 score: 0.707

⇒ *KNeighborsClassifier*(*n_neighbors*=10)

Mean precision: 0.785

Mean recall: 0.763

Mean f1 score: 0.764

⇒ *SVC*(*kernel*='rbf', *decision_function_shape*='ovo')

Mean precision: 0.443

Mean recall: 0.364

Mean f1 score: 0.314

As the results show, the random forest algorithm is the one that achieves the best performances, even with default settings, on both holdout and 10-fold crossvalidation. We thus decided to consider it as the main tool to build our solution.

5.2 Multiclass classifier

This is the first algorithm involved in our final model.

It somehow represents the core of our thesis, it models the starting idea: go deep inside bot identification and search and classify similar behaviours among them.

We started from the baselines identified earlier and we tried to refined them, adjusting to fit our need.

The workflow was the same as before, starting from the data with basis and handcrafted features, we took the best looking algorithm from the baselines pool and performed hyperparameters tuning on it, supported by a validation technique.

5.2.1 Dataset

During this phase, we used the previously described dataset 3.4 with its five different labels. The algorithm was fed with 26,357 samples and 38 features. the amount of records were light enough to consider K-fold crossvalidation, without slow the validation down too much.

5.2.2 Model

We found ourselves in the situation in which we had some brand new features and we didn't know how useful they were. Obviously, we could appeal to heathmaps or other tools, to highlight the correlations among variables and targets. However, the model we wanted to develop was the Random Forest, which proved to perform well with F1 score. Since this kind of model exploits its criteria to employ the features, we needed to prove them with a direct approach.

A useful advantage of the Random Forest algorithm is the ability to provide a feature ranking, according to its splitting criterion. We retrieved the 12 most important attributes, in order to see if we would have found some of the ones coming out from feature engineering. The algorithm ranking ranked the features this way: 1. *favourites_count* (0.199402), 2. *followers_count* (0.110830), 3. *statuses_count* (0.100810), 4. *avg_len* (0.058260), 5. *freq* (0.055419), 6. *friends_count* (0.043405), 7. *ret_perc* (0.039126), 8. *tweet_intradistance* (0.033718), 9. *max_ret* (0.030654), 10. *min_len* (0.029731), 11. *NSFW_words_score* (0.029323), 12. *default_profile* (0.022897).

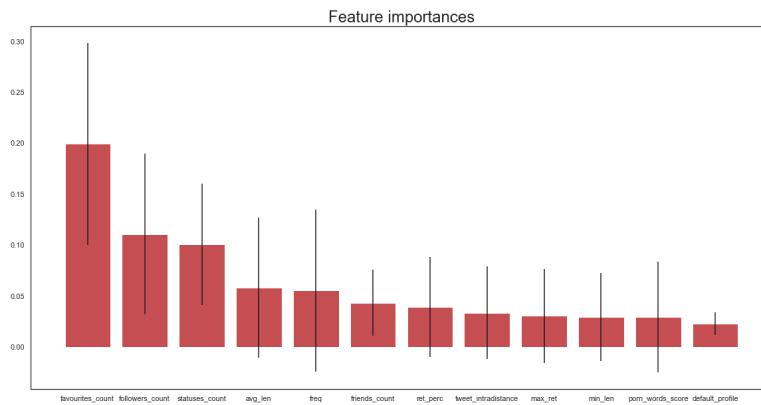


Figure 5.1: Random Forest feature ranking

As Figure 5.1 shows, we could find some of our crafted features inside

this list: lots of tweets descriptive features (*avg_len*, *freq*, *ret_perc*, etc...), as well as the *tweet_intradistance* attribute and the *NSFW_words_score*. This picture confirmed us that the idea behind those features was useful.

Since those attributes were thought to belong to different clusters, we decided to try several combinations of those feature clusters, validating the model on them with a crossvalidation. The purpose of this stage was to see if some groups of features were enough to describe the real problem, or if some group would shown up as irrelevant. To face this evaluation, we performed a light-weighted Grid Search, which is a method that takes desired ranges of hyperparameters and tries all the possible permutations of them, looking for the best combination, in terms of a certain metric.

We are talking about a light-weight version of this tool, because we just went through different numbers of tree estimators in the forest. The different feature groups are not considered as hyperparameters and are not handled by the Scikit-learn implementation of the Grid Search. We had to manage the different training by our own, looking how the test score would have changed along with the increasing number of estimators and the different set of features.

Grid Search uses crossvalidation to find the better estimators for the models, and this approach was right for our situation. Due to the multiclass nature and some imbalances with the labels, we decided to follow the F1 score metric to asses the value of our model.

The features were organized in clusters, as described in Chapter 4. We had the user features, the descriptive features, the intrinsic features, the extrinsic and the image features. Then we tried the model with the entire set of 38 attributes. As shown in Figure 5.2, the best configuration seems to involve the whole set of features, as it reaches these scores, with 70 estimators: *Precision* = 0.946, *Recall* = 0.945, ***F1*** = 0.943.

The model has been tested with the default value for the maximum depth in the trees, which is set to 'None'. It means that the trees are expanded until every leaf is pure, or all leaves contain one sample. Tuning the hyperparameters wasn't the goal of this experiment, since its aim was only to find the best configuration of features to fit the model with.

We then continued with a proper Grid Search over the whole number of features.

The Figures 5.3, 5.4 show how the average F1 score, measured on 10-fold crossvalidation, changes with the increasing of the number of estimators in the forest. The different coloured lines represent *max_depth* hyperparameter. The first Figure (5.3) shows the Grid Search results, with the *gini* splitting criterion. The second one (5.4) represent the situation having *entropy* as

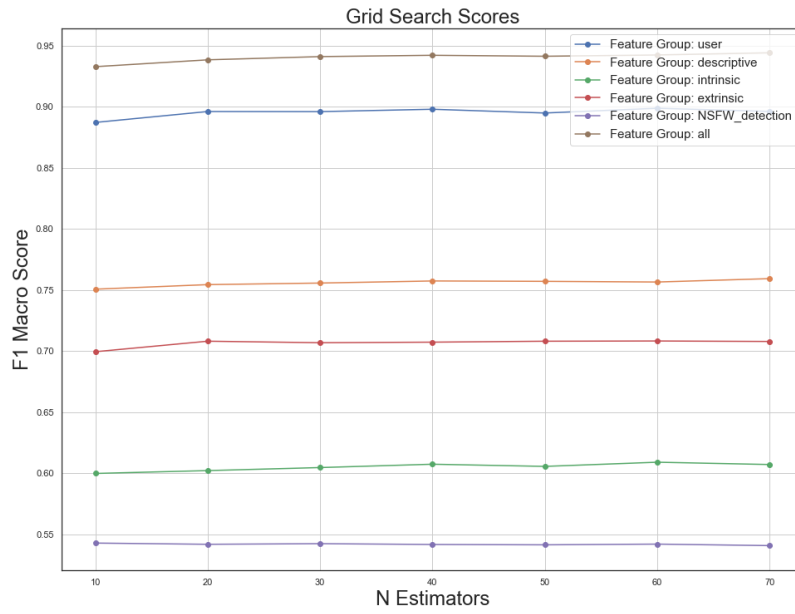


Figure 5.2: Performance over different feature clusters

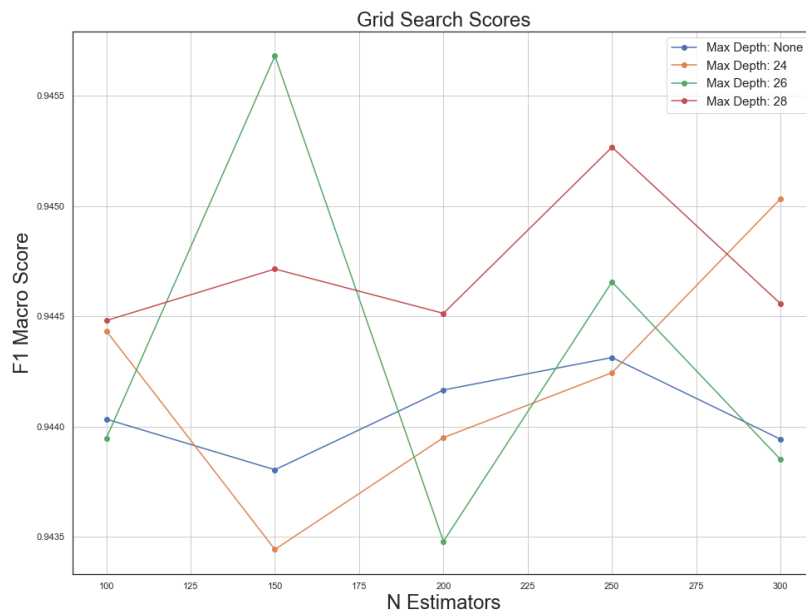


Figure 5.3: F1 scores with 'Gini' criterion

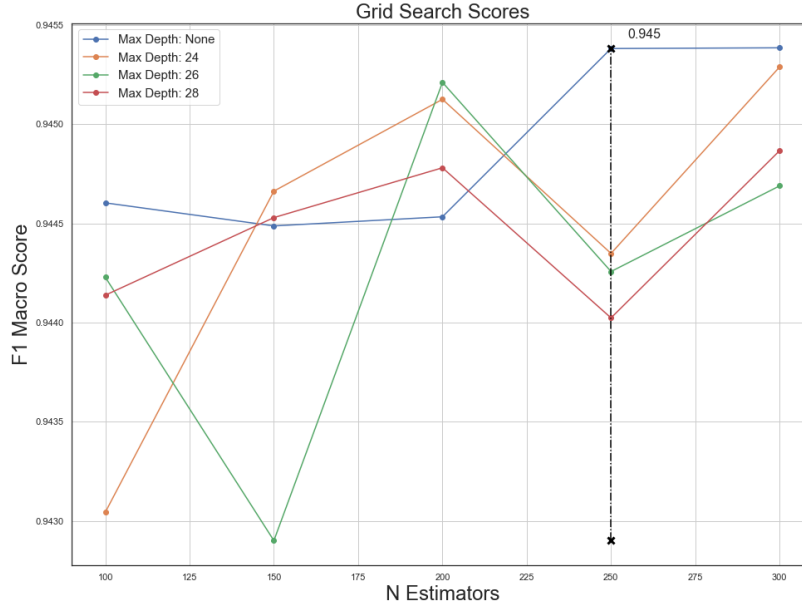


Figure 5.4: F1 scores with "Entropy" criterion

a splitting choice. We combined five numbers of estimators (100, 150, 200, 250, 300), together with four different depths of the trees (None, 24, 26, 28) and the two above-mentioned splitting criteria.

We could find a local maximum in 5.3, marked by the green line ($max_depth = 26$), matching $n_estimators = 100$. Observing the overall trend for that line, we couldn't rely on that peak, because it seemed to be a lucky hit, due the following lower values accomplished by this configuration. The entropy criterion seemed to be a more stable adjustments of scores, especially with the *None* setting for max_depth hyperparameter.

The highlighted point represent what we pick as our current configuration, and it sees 250 trees, the longest possible depth and the entropy mode for choosing the features on witch separate the trees.

We decided to stick with the following configuration without continuing with further explorations, because we could observe a flat improvements from 250 to 300 estimators, with that depth. Furthermore, we were aware that, in the generalization stage, the computation of the image features and the calls to the Twitter APIs would have been taking much more time than the models' predictions, so we decided to have a light model, in order to keep the prediction time fast enough, which is currently measured in 33ms.

The first algorithm of our solution was completed and ready to be combined with the following two models.

5.3 Binary Classifier

Since our dataset was pretty balanced and we couldn't retrieve much more genuine accounts, we didn't want our instrument to treat this category of users just as one the other bot kinds. It was important to perform a previous filter that was able to give importance to the separation between bots and genuine accounts.

We was inspired by the work made with Botometer[3], which involved a binary labelled dataset, with bot and genuine accounts. They built their features, grouped them in six main categories, then they ran a Random Forest algorithm per group.

We already had our feature engineering done, so we decided to test it on this new task.

In order to not to build a poorer version of our multiclass model, we didn't want to use a reduced copy of our dataset, stratifying it by stripping random bots from it. We needed a balanced dataset, with about the same amount of genuines and bots. So, we had to gather more human ids, because if we would have relied on the accounts we already had (3661), we would had disposed on a training set with about 7000 entries, that would have been too small to perform a relevant filter.

5.3.1 Dataset

The dataset we used for this classification was composed by part of our collected records and by some entries from the Carvelee-2011 dataset, which contains 22,223 content polluters and 19,276 legitimate users, both collected through a social honeypot, as described in their paper[4]. This dataset has been involved to build Botometer as well, but we decided to use it only partially, mixing it with our retrieved accounts.

We setted the APIs to retrieve up to 6,000 ids for both genuines and bots, from the Carvelee list. The process provided us 5,161 legitimate user ids, and 5,297 general bot ids (without inner classifications), because some accounts have been deleted in time. Then we added some new records, randomly sampling our data. This led us to reach 7,660 human accounts and 7,795 bots, forming a new dataset of 15,455 entries, with binary target.

The feature vector we used is the same that came out from the feature engineering process^{4.7}, except for the specific characterizing features, that

weren't considered, because crafted for the inner separation among bots. We excluded the image features (*NSFW_profile* and *NSFW_avg*) and the extrinsic keywords scores (*NSFW_words_score*, *news_spreaders_words_score*, *spam_bots_words_score*, *fake_followers_words_score*, *genuine_words_score*). It has been fitted with 32 features.

5.3.2 Model

Since the different purpose of this model, we couldn't rely on the same baselines tested with the multiclass dataset. We wanted to evaluate new raw algorithms for this binary classification. Another round of crossvalidation, with default settings, was performed on this new dataset, exploiting the binary nature to show the *Area Under the Curve* score (AUC). Area Under the Curve represents the goodness of a classifier, in terms of the integral of the *Receiver Operating Characteristic* (ROC curve), defined over the variation of a decision threshold. The motivation behind the adding of this new metric is that we had a balanced binary dataset, and this metric is a good fit for this kind of problem. Moreover, Botmoter claims to have accomplished an AUC of 0.95, on a 10-fold crossvalidation test.

The ROC curve lies in a bi-dimensional space, which has the *True Positive Ratio* ($TPR = \frac{TP}{TP+FN}$) on the Y-axis, and the *False Positive Ratio* ($FPR = \frac{FP}{FP+TN}$) on the X-axis.

We wanted a term of comparison, so we evaluated this metric even in the baseline stage. Here there are the results of this process:

⇒ *RandomForestClassifier*(*n_estimators* = 10, *criterion* = 'entropy')

Mean AUC: 0.916

Mean precision: 0.879

Mean recall: 0.793

Mean f1 score: 0.824

⇒ *LogisticRegression*(*fit_intercept*=True, *max_iter*=100, *penalty*='l2')

Mean AUC: 0.792

Mean precision: 0.694

Mean recall: 0.759

Mean f1 score: 0.723

⇒ *KNeighborsClassifier*(*n_neighbors*=10)

Mean AUC: 0.835

Mean precision: 0.779

Mean recall: 0.750

Mean f1 score: 0.760

✎ `SVC(kernel='rbf', decision_function_shape='ovo')`

Mean AUC: 0.583

Mean precision: 0.620

Mean recall: 0.364

Mean f1 score: 0.095

Once again, Random Forest won the comparison with the other baselines. We decided to let it perform this job and try to improve its AUC score, keeping one eye on the overall score too.

The algorithm has had its parameters tuned during the validation phase. We decided to stick with 10-fold crossvalidation, as it was done for the baselines.

After several Grid Search runs, the last round computed had this hyper-parameters to combine together:

✎ `n_estimators = [100, 115, 130, 150, 175, 200]`

✎ `max_depth = [None, 26, 28]`

✎ `criterion = 'entropy'`

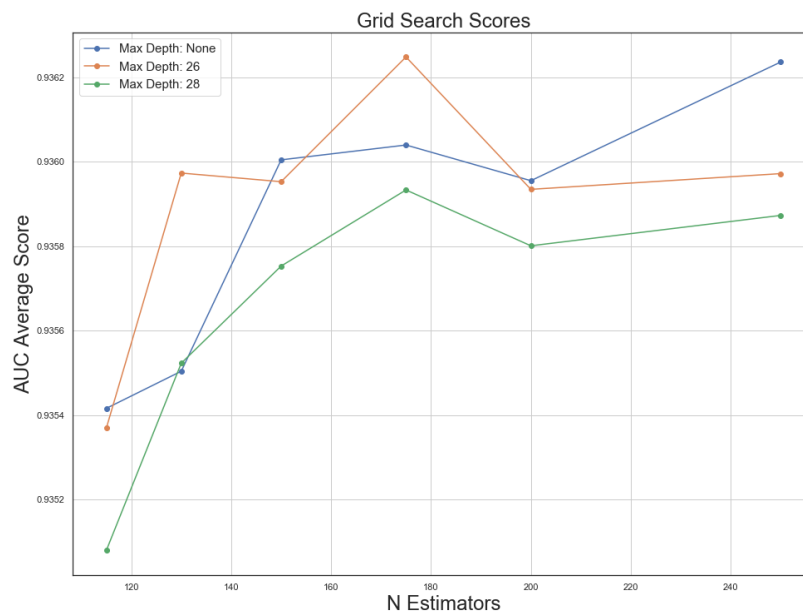


Figure 5.5: Grid search results

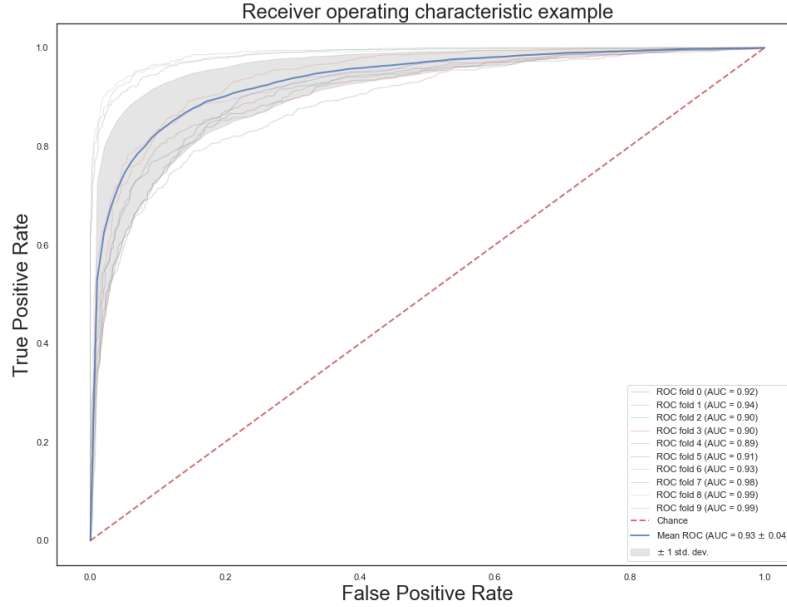


Figure 5.6: ROC curve

As we can see in Figure 5.5, the AUC is increasing with the number of the estimators in the forest. We decided to stop at 250, which corresponds to the highest AUC score, with a low risk to perform overfitting.

The AUC obtained with this arrangement is equal to 0.936, as shown in Figure 5.6, which is a positive accomplishment, considering that it will be used only as support for the identification of humans among bots, but we didn't craft specific features as the ones involved in the Botmoter project and we didn't have the same amount of data neither.

The final model has been fitted with the whole data, with these settings: $n_estimators = 250$, $max_depth = None$ and $criterion = 'entropy'$. The max_depth parameter states how deep the tree should be expanded. When set to *None*, the algorithm tries to expand them until every leaf is pure. We were looking for this approach, because we wanted this filter classifier to ensure the hardest classifications possible, among a soft classification system. We wanted some loose detections, when there were no certainties. But we wanted the prediction to be as sure as possible, when no uncertainties were met.

Once we fitted the model, it was ready to be the second element entering in our solution pool.

5.4 Text classifier

5.5 Stacking meta-classifier

We had three models, each with different purposes, but they had to cooperate for the bots' behaviour identification. The initial idea was to use only the multiclass Random Forest to classify the bot categories, using the other two models as meta-models to build extra features with their outcome. Those features would have had the dataset enhanced, but their meaning would have been bounded to the multiclass classifier limits. We wanted to give the right importance to each model, hoping they would help each other to better distinguish the patterns and to better model the real problem.

We thought about several methods to exploit their strengths and combine them. The first idea was to build a pipeline with weights for each classifier. The binary Random Forest would have been the first filter between humans and bots. Following this lead, the multiclass and the text classifiers would have been weighted in order to assign the final label.

This approach would have been less empirical than other tools we could dispose. We decided to rely on different methods to blend the outcomes of our three algorithms. In particular, we thought about a genetic approach and a stacking ensemble with a meta-classifier. We wanted to evaluate the performance obtained by these methods and chose the one that fitted our need.

Both the genetic and the meta-model were trained with holdout technique, splitting the whole dataset into training and test sets. The 80% of the samples ended up into the training set, the 20% in the validation set. We had a training set for the ensemble models that contains 5022 entries. The data that fed the stacking methods were the predictions of the tree classifiers, over the validation set. In order to make those predictions without cheating, we couldn't use the models that were already fitted with the whole data. We had to retrain them with the 80% of the records. We didn't perform further Grid Search to find the best hyperparameters in this stage, because the final script that we were going to assemble was taking into account the entire dataset to train the models. Furthermore, this small variation, in terms of amount of training data, wouldn't have led us into a misinterpretation of the problem, if we had kept the same hyperparameters found earlier. We decided to stick with the configurations already found and to train the model with fewer data.

Once the models were fitted, we used the *predict_proba()* method of the Scikit-learn implementations of the classifiers, in order to retrieve "soft clas-

sifications”. We didn’t want our model to assign a strict label to an unseen sample, indeed, we were interested in the percentage of categories membership. The `predict_proba()` method computes the probability, for a sample, to belong to the highlighted target, by considering the impurity of the leaves inside the forest. We used this method to construct the output vectors needed to train the stacking models.

In order to combine the outputs of three classifiers properly, we had to homologate them. With the multiclass and the text models there weren’t issues, since their outcomes are stored into vectors with five elements, containing probability values, one for each category. The case to handle was the binary one, because the output of that Random Forest was represented by a 2-sized vector (the bot probability and the genuine one). We extended the array adding three more cells. Then we took the bot probability and spread it on the four cells that didn’t mark the genuine probability. The reason behind that is that our binary classifier didn’t take care about the inner separation in bot behaviours. It was uncertain about the belonging category, it just used to recognize bots, so each non-human category should had been assigned the same coefficient, and their sum must had match the original probability to be detected as a bot. For instance, if the binary classifies detects an account with 80% chance to be a bot and 20% to be genuine, it will result in this final outcome vector:

Bot	Genuine
0.8	0.2

Original binary output

NSFW	News-Spreader	Spam-Bot	Fake-Follower	Genuine
0.2	0.2	0.2	0.2	0.2

Adapted multiclass output

Each sample of this new dataset contains 15 elements, 5 soft predictions (one for each category) for each classifier (binary, multiclass, text-based).

New sample

Binary probability					Text-based probability					Multiclass probability				
P0	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12	P13	P14

A new training set was born and it was built with the soft classifications of the models in the pool, over the validation set. It was ready to proceed and to serve the ensemble models.

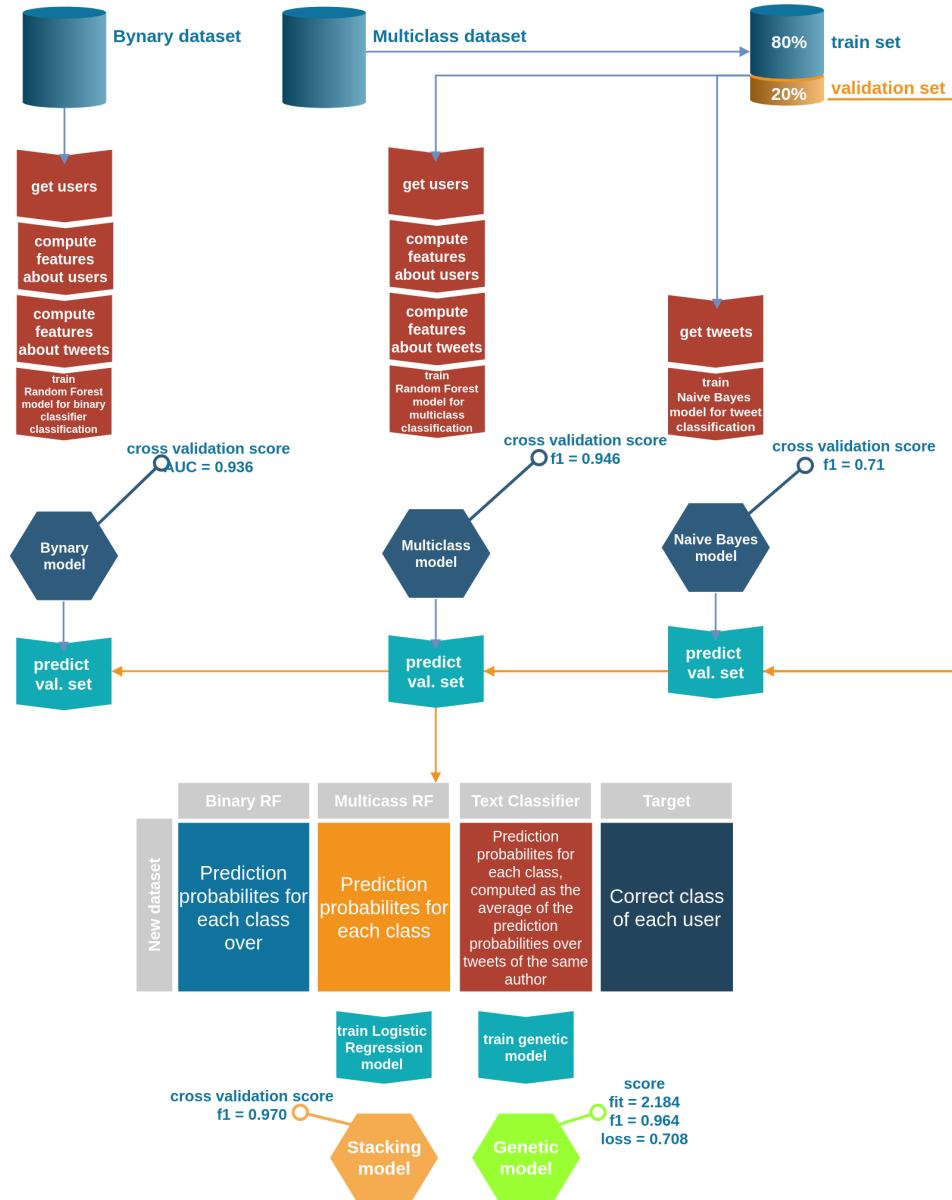


Figure 5.7: Training pipeline of theNodv stacking ensemble

The pipeline for perform the training of the ensemble models is resumed in Figure 5.7.

After several parallel attempts were done, the blending system was built

with the meta-classifier.

5.5.1 Genetic algorithm

This approach started as a side way, when we were already testing the stacking ensemble.

The idea behind genetic programming, is to emulate the natural species evolution, by encoding the the *chromosomes* in the process with data structures. The chromosomes represent the possible solutions for the problem and they have to "evolve", in order to get fitter and fitter for the goal. Several operators must be determined to perform this evolution. Once a first *generation* of feasible chromosomes has been formed, they have to be evaluate according to a *fitness function*, which asses how well a chromosome faces the problem. The best portion of chromosomes are picked to be part of the next generation, and this is called *elitism*. The solutions left are given a probabilities to join the elite ones, in order to form a new generation with about the same size as the previous. This step is called *selection*. The chromosomes picked in the selection stage are assigned a high *crossover* probability. The crossover operator handles the "born" of new chromosomes, mixing parents alleles in a certain way. The mixing method is highly correlated to the chosen encoding strategy. Each newborn is give a low probability to undergo a mutation. This step often seems useless, but it's pretty important, in order to explore a higher spectrum of solutions, which couldn't be expanded by the mating operators only. After the new population has been accepted, it is ready to be validated through the previously define fitness function. The loop holds, until a solution is found, or, like in our case, the process sticks to a local or global maximum.

Genetic operators

We setted the genetic algorithm with the support of the Deap library for Python, setting these operators:

- ⇒ *Encoding*: each chromosomes represented a weighting vector for the outcome of our three classifiers. Each allele of the chromosome was float valued, with numbers between 0 and 5, generated randomly, with a uniform distribution. We started with normalized weights, but the spectrum of the solution explored was way too poor to fit the needs. This range was given after observing the weights that the Logistic Regression model were assigning to the inputs received, that was wider

and involved even negative values. We randomly generated 200 chromosomes for the initial population, with this form.

Chromosome														
Binary Weights					Text-based Weights					Multiclass Weights				
w ₀	w ₁	w ₂	w ₃	w ₄	w ₅	w ₆	w ₇	w ₈	w ₉	w ₁₀	w ₁₁	w ₁₂	w ₁₃	w ₁₄

⇒ *Fitness evaluation*: the fitness function that assessed the value of the solutions was somehow similar to the one used in the other stacking method. We applied the weights of our chromosomes to the samples in our dataset.

For each sample, we made pairwise additions, among the outputs of different classifiers, multiplied by the chromosome's weights, for the same category:

$$\begin{array}{c}
 \begin{array}{c} \text{Binary Components} \\ \hline p_0 * w_0 \dots p_4 * w_4 \\ \hline \end{array} \\
 + \\
 \begin{array}{c} \text{Text-based Components} \\ \hline p_0 * w_0 \dots p_4 * w_4 \\ \hline \end{array} \\
 + \\
 \begin{array}{c} \text{Multiclass Components} \\ \hline p_0 * w_0 \dots p_4 * w_4 \\ \hline \end{array} \\
 = \\
 \begin{array}{c} \text{Resulting Prediction} \\ \hline \mathbf{p_0 \ p_1 \ p_2 \ p_3 \ p_4} \\ \hline \end{array}
 \end{array}$$

In order to stick to the probabilities nature, the computed prediction had been normalized.

That prediction has been compared with the known real target for the examined sample. Since the targets of our dataset aren't soft valued, we took the maximum probability of the computed prediction

to make the comparison with the actual class. Our fitness function aims to favourite those solutions which maximizes the F1 macro score, as it has been for the validation of the classifiers, until this stage.

During this process, the problem we had to face was that we wanted to produce soft classifications, because we knew that our collected data presents similar patterns within the same categories. This means that the algorithms easily classify our test set, because of the distinctive traits found for each target. In order to mitigate the real test error, over unseen samples, we wanted the prediction to be as smooth as possible, without confusing the F1 score interpretation. In addition, the multiclass Random Forest was setted with the maximum depth possible for its estimators, that implies a drastic reduction of impure leaves, leading to hard classifications.

We faced the problem involving a smoothing factor to our fitness function.

When computing a sample, we populated a Confusion Matrix of the prediction, using the above-mentioned method to match predicted and actual classes. The matrix helped us computing the F1 macro score easily. At the same time, we counted every hard classification, marking as 'hard' every computed prediction that contained a probability greater or equal to **0.8**, among its five stored values. This count was used as a penalty, it has been averaged for the number of samples, and then subtracted to the computed F1 score of that chromosome. In order to privilege the maximization of the F1 factor, instead of the minimization of the penalty, the final fitness function assigned this score to each chromosome:

$$Fitness = 3 \times F1_score - Penalty$$

This way to operate didn't affect the overall F1 of the sample, since penalizing hard classifications didn't discourage the system to look for values high enough to have a dominant category in the prediction.

Once every chromosome has been evaluated, they could proceed to the next steps of the algorithm.

- ➡ *Selection:* the selection phase handles the choice over which chromosomes pick for mating. Several pre-implemented methods are available, but we used the tournament method. It works selecting the size K of the tournament, which we chose to be 3. Then, it randomly selects K (3) chromosomes from the population and places it inside a pool.

Then it compares their fitness. The chromosome with the best fitness has probability p (the crossover probability) to be selected for mating. The second has $p*(1-p)$ chance to get selected, the third $p*((1-p)^2)$.

⇒ *Crossover*: The crossover probability has been setted to 95%. The crossover operator wasn't something already implemented by the library, as for the fitness function. Our operator used to produce two brand new chromosomes for the next generations. The first child is the unweighted mean of its parents:

$$[x_0, x_1, \dots, x_{14}] \oplus [y_0, y_1, \dots, y_{14}] = \left[\frac{x_0 + y_0}{2}, \frac{x_1 + y_1}{2}, \dots, \frac{x_{14} + y_{14}}{2} \right]$$

The second child is the weighted mean of its parents, computing the weights with respect to the fitness of the two mating chromosomes:

$$f_x = \frac{fitness_x}{fitness_x + fitness_y}$$

$$f_y = \frac{fitness_y}{fitness_x + fitness_y}$$

$$[x_0, x_1, \dots, x_{14}] \oplus [y_0, y_1, \dots, y_{14}] = [x_0 * f_x + y_0 * f_y, \dots, x_{14} * f_x + y_{14} * f_y]$$

The retrieved children used to be part of the upcoming generation.

⇒ *Elitism*: This part was necessary, in order to not lose the best solutions found so far. It is a sort of insurance, which guarantees to keep, at least, the best situation until this stage, and to let it take part of the next generations of solution. We preserved our three best chromosomes for each generations and move them to the next stages.

⇒ *Mutation*: The mutation probability is generally setted to low values, like what happens in nature. It represent the error in DNA replications from the parents and it shouldn't reach the 1% of probability to occur. Although, we wanted to force some mutation, because, as said before, we needed a wider space of solutions and the elitism helped us in containing the damages of such mutations. In the worst cases, all the chromosomes have had been damaged and resulted as useless, but the elitism had preserved the best ones and kept it untouched. So we imposed a 45% of mutation probability, for each newborn solutions, before entering the pool.

Our mutation operator was a decoration of the value changing method already implemented: we randomly used to pick three elements from the chromosome and set them to zero.

Results

After several runs of the genetic program, with boosted starts (the best solutions found at the previous run were placed inside the first generations of the following runs), we stuck in a maximum of the F1 score. In the last run, from the 5th generation there was no improvements in the fitness of the best solution. We selected the fittest chromosome, whose scores were:

⇒ *Weights:*

Binary Weights				
0.0	0.234	0.0	4.145	2.668
Text-based Weights				
4.867	2.808	4.680	2.800	2.325
Multiclass Weights				
4.462	2.129	2.162	2.317	0.163

⇒ *Fitness:* 2.184

⇒ *F1 score:* 0.964

⇒ *Percentage of hard classifications:* 70.8%

The solution found didn't take into account the first and the third element of the input, which correspond with the outcome of the binary classifier, for both the NSFW and the Spam-Bot categories. It seemed legit, because the main difference between that model and the multiclass classifier, except the binary task, is the lacking of the NSFW specific features. It makes sense that this model wasn't able to contribute to the separation between those kinds of bot. That feature vector should had been involved in the final model, simply as a weight for the weighted mean of the probability of each classifier, for each sample to predict over.

Although the result seemed promising, we knew that a "simple" weighted mean of the outcomes of the classifiers weren't enough to describe the problem. The risk to have performed overfitting over our data was high.

We decided to test this final weighting vector with new unseen samples, selected manually for the purpose. The comparison with picked samples didn't match the expectations given by the metric evaluation, so we tried a different and more sophisticated stacking method.

5.5.2 Logistic Regression

The reason behind the choice of a meta-classifier is that we wanted a more complex way to perform inner weighting of the outcomes that we had from other models. A simple weighted mean wasn't enough for this purpose. Furthermore, implementing a logistic-like loss to evaluate the fitness of a genetic algorithm would have meant to apply the Logistic Regression training model, without performing gradient descent, but with a genetic approach. It would have been just unnecessary and computationally expensive. Thus, we discarded the idea of using the genetic programming to emulate a Logistic model, even if the smoothing factor used for that try was a good insight for our task. In order to mitigate the lacking of soft classifications, we chose to rely on the regularization factors that belong to the training algorithm of the Logistic Regression. This kind of models is often involved in stacking other classifiers, with a binary purpose.

Dataset

The same training set has been used for train both the Genetic and Logistic models. Since we were managing a multilabel datasets, we knew that the ensemble meta-model would had been adapted to this job. The most common tool used for stacking purposes is the Logistic Regression, which performs well on binary separations. We decided to test this model on a multinomial approach, with a softmax activation function, instead of trying the already visited One-vs-Rest method.

Comparison with Random Forest

We didn't wanted to blindly select this model over some others tool, especially over Random Forest, which proved us to perform well in multiclass classifications. Thus, we tested these two algorithms with the new dataset. We ran some default configurations of the models, in order to have a raw comparison to trace a line between them.

Figure 5.8 shows the early convergence of the Logistic model's F1 score, with low maximum iterations. The model has been tested with Lasso penalty and two different solvers, but the results, over the increasing of the training epochs, are way similar. The Random Forest, as shown in Figure 5.9, needs lot of estimators to top the performance of the Logistic Regression. It was tested over the number of features to consider for the splits ($\log(|features|)$ and $\sqrt{|features|}$), along with an increasing number of trees.

We preferred the lighter Logistic model over the Forest ensemble, for

htp!

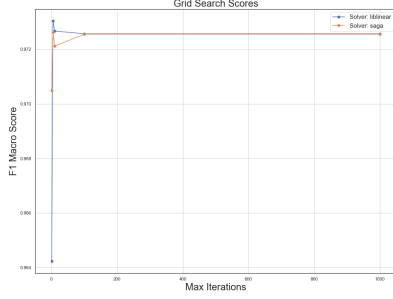


Figure 5.8: LogReg with raw settings

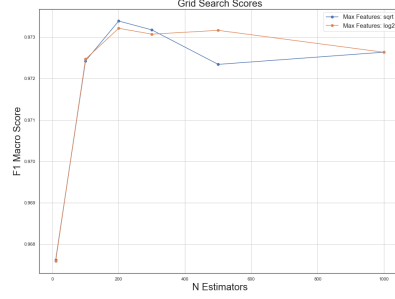


Figure 5.9: Random Forest with raw settings

Figure 5.10: Stacking models comparison

performance matters and prediction time: Logistic predictions take 0.6ms in average, in comparison with Random Forest with those settings, which averages 100ms to provide an outcome over a single sample.

Hyperparameters tuning

We tried two regularization terms for the Logistic model and several numbers of maximum iterations for the training algorithms. The regularization terms are parameters computed in addition with the minimization of the characteristic loss function. Their purpose is to avoid the weights to explode and the model to become more sensitive to noisy data. In other words, they are involved to prevent overfitting. The idea is that the loss function, gets modified as follows

$$\mathbf{L}(\mathbf{w}) = L_D(w) + \lambda L_W(w)$$

Where $L_D(w)$ represent the error on the data and $L_W(w)$ is the term representing the model complexity. In general, smoother weights implies lower model complexity. The lighter the complexity, the lower the variance of the model and the risk perform overfitting. The parameter λ has to be tuned with a validation method.

The penalties that we explored were:

📦 *Lasso* (L_1):

$$\mathbf{L}_1(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|_1$$

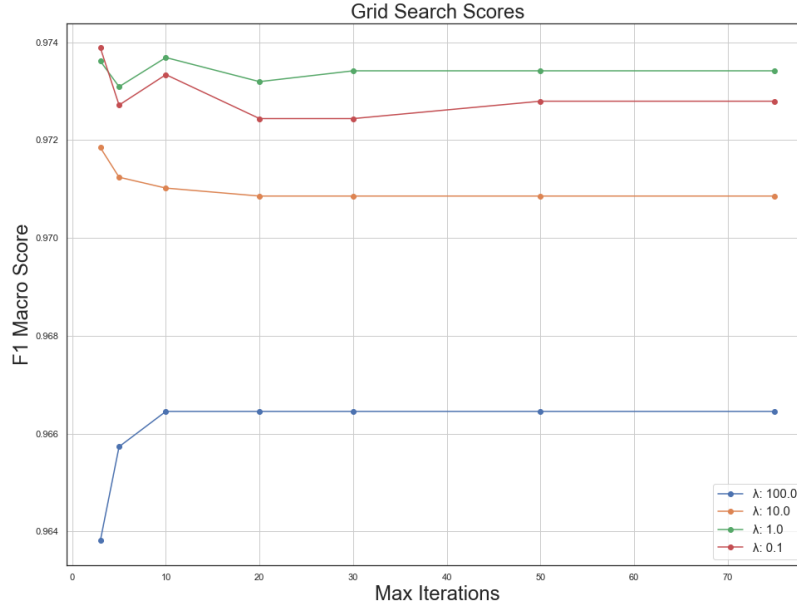


Figure 5.11: Lasso, $\lambda = [0.1, 1, 10, 100]$

where $\|\mathbf{w}\|_1 = \sum_{i=1}^N |w_i|$

This regularization function is non-linear and doesn't provide a closed-form solution. It tends to cut out some features from the model, yielding to sparse and lighter model. It can be seen as an implicit way to apply features selection.

🔗 Ridge (L_2):

$$\mathbf{L}_2(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

where $\|\mathbf{w}\|_2^2 = \sum_{i=1}^N w_i^2$

This softer term tends to shrink the weights, keeping the loss function quadratic in \mathbf{w} and closed form solution exists.

Figure 5.11 highlights the slightly better results obtained with the Lasso penalty, with unitary λ coefficient (Lasso F1 score: 0.973). The ridge penalty needs to be weakened ($\lambda = 0.1$) in order to top the Lasso performance, which is a compromise hard to deal with. The smaller the regularization coefficient, the higher the model complexity, as said before. Moreover, we decided to gather further consideration, by looking inside the weighting applied by those two terms.

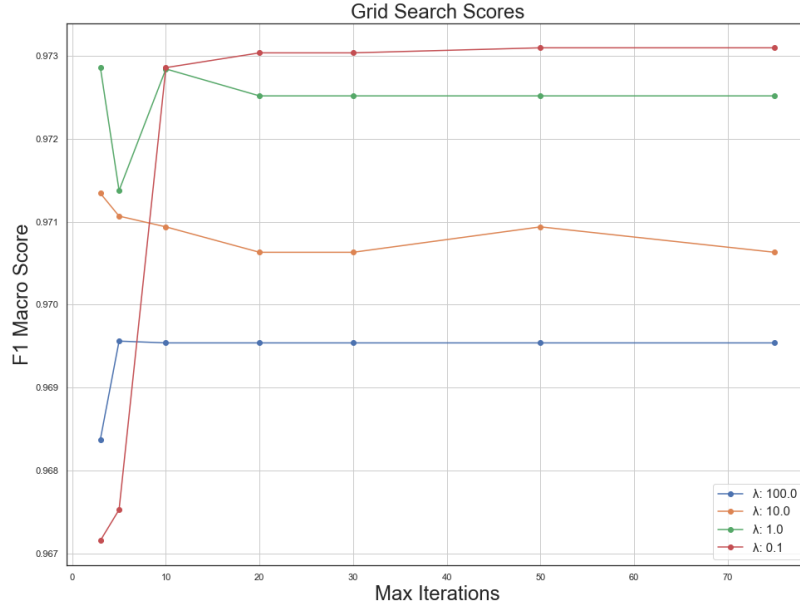


Figure 5.12: Ridge, $\lambda = [0.1, 1, 10, 100]$

Ridge regularization				
Binary Weights				
$5.1e^{-16}$	$5.1e^{-16}$	$5.1e^{-16}$	$5.1e^{-16}$	$-8.8e^{-17}$
Text-based Weights				
$-2.2e^{-15}$	$-1.4e^{-14}$	$-5.1e^{-15}$	$-4.2e^{-15}$	$1.1e^{-16}$
Multiclass Weights				
$-3.0e^{-15}$	$2.0e^{-15}$	$-1.5e^{-16}$	$5.7e^{-15}$	$-4.2e^{-15}$

The Ridge regularization leads to very small weights, and negative ones too. Even with unitary λ coefficient, it is hard to distinguish a discrimination among features. This approach would have yielded a smoother model, but with the ability to give a chance to every classifier to distinguish among targets. We wanted to get some further insights from the other weighting.

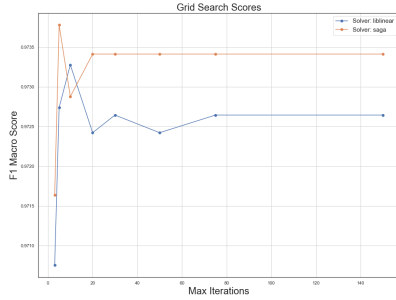


Figure 5.13: Up to 150 iterations

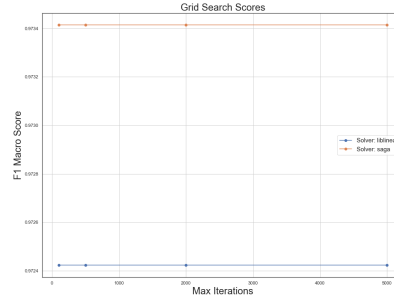


Figure 5.14: Up to 5000 iterations

Figure 5.15: Lasso Logistic Regression scores over solvers

Lasso regularization				
Binary Weights				
0.0	0.0	0.0	0.0	1.391
Text-based Weights				
0.020	1.297	0.0	0.182	0.952
Multiclass Weights				
1.316	0.480	1.457	2.087	0.435

The L_1 term, as expected cut out some features from the model. Looking at the excluded attributes, we noticed that the regularization caught the nature of the binary classifier. It wasn't supposed to give some contributes to the bots inner separation, its main purpose was to detect humans from content polluters. Lasso seemed to understand this behaviour and decided to not consider the opinion of that classifier, when it comes to bot categories. Another insight got from L_1 is about the treatment of the text-based classifier. the Naive Bayes model couldn't distinguish with certainty between NSFW and Spam-Bots, since they act in a similar way. They just tweets click-baiting links, with catchy captions. A "blind" classifier struggles in understand the nature of those links. This is the reason that almost two contributes of Naive Bayes had been discarded from the stacking model.

Since we knew our dataset and we were aware of the bias it might contain, we preferred a lighter and sparser model, over a more complete one, even when the F1 scores used to match. We wanted our model to infer on new unseen data and to be ready to give a representative statistical description of the actual situation on Twitter. We had to be far-sighted and not to recline

on the accomplishments of the 10-fold crossvalidation. We thought that the Lasso model would have been performing better in out-of-box predictions.

We kept the L_1 penalty, with $\lambda = 1$ and proceeded with the hyperparameters tuning.

Figure 5.13 shows the trend in the F1 score, along with the increasing number of iterations, applying the *SAGA*[1] solver (a variant of the *Stochastic Average Gradient*[5] optimization that supports Lasso penalty) and the *LIBLINEAR*[7], an open source library for large-scale linear classification.

As it can be seen in Figures 5.14, by increasing the number of maximum iterations, up to 5000, the performances remain stable with every solver. The algorithms seem to not improve after 75 maximum iterations settled. Moreover, the SAGA solver gains slightly better results, in terms of F1 score, as it reaches 0.9734 in this metric, compared with the score obtained with LIBLINEAR solver, which is 0.9727.

The final Logistic Regression meta-classifier has been fitted with the Scikit-learn library, with this setting:

```
LogisticRegression(max_iter = 100, penalty = "l1", solver = "saga", C = 1, multi_class = "multinomial", fit_intercept = True, class_weight = "balanced")
```

The C parameter stands for $\frac{1}{\lambda}$, the regularization coefficient. The balanced class weights means that the model automatically adjusts class weights inversely proportional to class frequencies in training set.

This setting obtained the following scores in a 10-fold crossvalidation:

⇒ *Precision*: 0.972

⇒ *Recall*: 0.974

⇒ *F1 score*: 0.973

5.6 Prediction pipeline

The final model is represented by a stacking ensemble with a defined execution pipeline for predictions.

As described by figure 5.16, In order to perform a prediction over a new sample the process is the following:

⇒ User and tweets data retrieving with Twitter APIs

⇒ Prepare data and perform features engineering to output multiclass probability prediction (**Multiclass Random Forest**)

- ⇒ Prepare data and strip multiclass features to output binary probability prediction (**Binary Random Forest**)
- ⇒ Prepare and treat text to perform text-based probability prediction (**Text-based Naive Bayes**)
- ⇒ Build new features vector with the stacked outcomes of the previous classifiers
- ⇒ Compute the final mutliclass probabilistic prediction with the meta-classifier (**Multinomial Logistic Regression**)

The produced probability vector assesses the nature of the examined user. It will be handled by a web application, in order to provide a useful classification tool for every internet user. The engine of this web application is mainly composed by a python script, which, given a Twitter user name, resumes this prediction pipeline and executes it, providing the classification.

This last lines anticipate the content of the following chapter of our thesis.

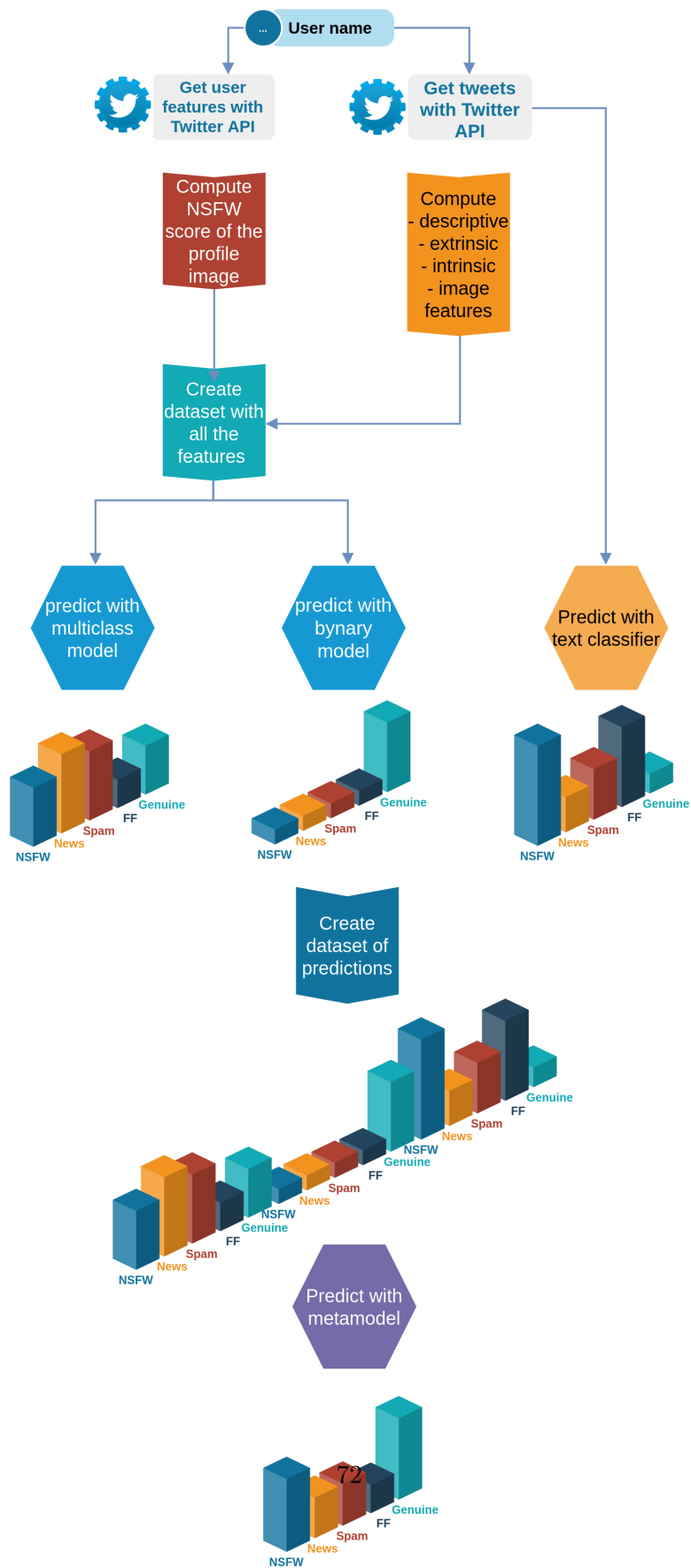


Figure 5.16: Final prediction pipeline

Chapter 6

Web application - BotBuster

6.1 Engine

6.2 Backend

6.3 Frontend

6.4 Deployment platform

6.5 Validation

Chapter 7

Implementation and Evaluation

Yes, you got it: finally, let's talk technology! If you are an attentive reader, you will have noticed that so far I restrained from talking about technology and implementation stuff. And that was intentional: doing a thesis is first and foremost a *conceptual* effort, meaning an effort that should require a lot of brainwork, thinking, reasoning, discussing, drawing sketches of ideas, constructing tables for making informed choices, and so on.

And you know what? If that is well done and well described, your reader, even if he/she is not tech-savvy or an expert in your topic, will understand you and be able to follow your reasoning and agree/disagree with the choices you propose. If instead you start too early talking about technologies, programming languages, protocols, fancy frameworks that your reader does not know and, even worse, explain your solution in function of these technologies, you will lose the attention of your reader. And there is nothing as bad as that.

Once you lose the attention of your reader due to too much geek talk, you will not be able to get the attention back. The consequence is that, even if you did the best project ever and come up with Nobel Prize worthy findings, your reader (perhaps your reviewer) will not notice, and you will not get the credit you actually deserve.

The lesson learned is: *defer* the tech talk as long as possible (too early = too dangerous), *single it out* from the rest of the work (so that who is not interested in the low-level details can skip it), and make it *self-contained* (so that who instead wants to read it gets all the details necessary).

7.1 Implementation

Here you can describe the technologies you use, put code example, describe all the details you feel are needed to enable the reader (with the necessary tech background) to understand. The goal of this section should be to enable your reader to re-implement what you did, perhaps with different technologies.

- ☐ Describe the *technologies* you use in your solution.
- ☐ *Motivate* possible technology choices.
- ☐ Copy and paste here the *architecture figure* you should already have included in Section ?? and extend it with the technologies you use for each of the modules.
- ☐ Provide insight into the most important *implementation problems* and how you solve them.
- ☐ If available, provide a link to an *online repository* holding the code of your prototype (ideally released as open-source software on GitHub or the like).
- ☐ Maybe you also want to share here some *UML diagrams* you drew before starting with the coding of the software.
- ☐ Provide evidence that your prototype *works*, e.g., screen shots, produced outputs, or similar.

7.2 Evaluation

This is a section that may be missing in a tesina, while for a thesis it is of fundamental importance. Even more: in some projects, the evaluation may even be a major contribution of the work and deserve an own chapter. If this is your case, then do so. For example, if you do an elaborated user study that requires careful literature study, design, planning, execution, data collection, data analysis, then you may want to make this effort also evident in the structure of the thesis by giving it an own chapter (remember that the structure of the thesis should already tell the reader a story).

7.2.1 Design of Evaluation

Explain here how you evaluate your solution, e.g., you do a controlled performance study in the lab using a cluster of 50 computers in a network, or

you do a simulation of an algorithm for which you first do some probing of some environment to fine-tune some parameters of the algorithm to have the simulation represent as real as possible situations, or you may do a user study, or... Here some options:

- ☐ *Theorem proofing*: if your work is of pure theoretical nature, you may want to accompany your theorems and corollaries with suitable proofs. Doing so requires good mathematical and/or algebraic skills.
- ☐ *Data analysis*: if you work on a topic that is related to Data Science, likely you will have a lot of data to analyze. Explain which data you are considering, how it is collected and prepared for the analysis, which kind of statistical analyzes you intend to use, why, etc.
- ☐ *Performance test*: if instead you develop a software prototype and claim that it works better than some exiting algorithm/software, explain which is your baseline to back your claim, tell how you want to compare your solution to the existing ones, which results you consider a success and which instead represent a failure, etc.
- ☐ *User study*: if your work involves real users in the evaluation of your work, explain how you select the participants, if they have to sign a consent form or not, if they need to obtain some form of prior training, which data you collect, how you guarantee their privacy and the security of the collected data, how you analyze the data, etc.
- ☐ *Simulation*: if you are not able to run your solution in a real environment and instead have to fall back to a simulation, explain how you set up the simulation environment, which assumptions you make, how you configure the simulation environment so that it resembles real situations, which exact data you collect, how you analyze it, etc.
- ☐ *Case study*: if the nature of your work does not allow a systematic data collection to back your claims, perhaps you want to elaborate on a case study that showcases the use of your solution in a real or fictitious application scenario. Explain the requirements of the case study, tell how realistic the case study is, show how your solution helps.
- ☐ ...

7.2.2 Metrics

Remember when I talked about the requirements and that ideally it would be good if the reader in the end of the thesis was able to use the list of

requirements as a checklist and to tick boxes? Well, this is where the reader should get the necessary tools to tick the boxes. Most likely, some of the requirements, claims and evaluation designs will need some specific metric to be able to tell if a requirement is satisfied or not. For example, you may want to measure response time for a time-critical service, or precision/recall for works on information retrieval, or individual quality attributes in crowd-sourcing, or...

- ☐ Define all the *metrics* needed by your evaluation designs.
- ☐ Tell how to *assess* the requirements and claims of your thesis.

7.2.3 Results

In this subsection you report on the results of the experiments/evaluation you perform. Report on all the important numbers for each of the metrics, on possible issues with running the code, etc. This is however not yet the place where to go into lengthy considerations on the meaning of values, this is for the next subsection. It's good to explain comparative results (A better than B in condition X, while in condition Y B is better), outliers (in one very specific situation A has an extraordinarily low/high performance), general statistics.

7.2.4 Discussion

Finally, here you discuss your results. That is, you discuss the *meaning* and *impact* of your result for the goals of your thesis. In other words, you *interpret* the results in light of your goals, expectations, intuitions, hypotheses. Did the prototype meet the expected performance? Is the achieved statistical significance reached to draw conclusions you would not be afraid of defending in front of a commission? Was the problem solved? Too slow? Too fast? Give the reader a feeling (as well as convincing arguments and numbers) for why you think some requirements are met while others may be missed.

Chapter 8

Conclusion and Future Work

So far so good. We are almost done. What is left is, well, just one of the most important chapters of the whole thesis, i.e., the conclusion. The purpose of this section is not to “conclude” the thesis in the sense to “stop” here. It’s rather to draw conclusions, that is, tell how well your work actually meets the requirements identified, answers the research questions, advances the state of the art. As such, this is perhaps the most important section! It may seem easy to just summarize a bit what you did and tell again what your objectives were when starting the work. But be aware that this can be much more difficult than it sounds, and you can expect your supervisor iterating with you several times over this same chapter. It is important that you show again your personal and professional maturity and your understanding of the topic. As you will see, some healthy self-criticism too is needed to make this chapter good.

8.1 Summary and Lessons Learned

Summarize here your work in about one page.

- ☐ Start from the initial *problem statement* or *research questions*.
- ☐ Summarize your *approach* and *methodology*.
- ☐ Recap the *lessons learned*.

8.2 Outputs and Contributions

Provide an overview of the outputs your project/work produced and then state what you think are the (research) contributions that advance the state of the art.

- ☐ List all the concrete *outputs* you produced (remember the discussion in Section 1.4).
- ☐ Copy/paste here the *list of contributions* you already anticipated in Section 1.4 (attention: outputs and contributions are two different lists; don't mix them).
- ☐ For each of the contributions, provide suitable *evidence*, drawing from the body of your thesis. For instance, if you claim that you did a formal proof of something, provide the exact number or name of the proof. If you promised subjective evidence for something, link this claim to the user studies you did. Etc. One or two sentences are enough for each of the contributions.

8.3 Limitations

This is where your self-criticism is needed. By now, I am confident you did a great work with your project and the writing of your thesis. So, compliments for that! You're almost done. But let's be frank: the work is not perfect. It simply cannot be, it never is. If it is, then not only I but also the whole commission of your defense will give you a standing ovation (I really would like to see this once). But in general there are just so many aspects of a research/thesis project that one would have to control or test, and with the limited time and resources available for these kinds of final projects it is just not possible to do everything.

In this section, you therefore tell the reader which aspects of your work may limit the impact or generalizability of your findings or contributions. As said, be frank. If you tell that you did a user study with only 10 people instead of 30 (which would make the findings stronger), you don't risk to give the impression you didn't do it well enough. Actually the opposite is true: if you don't tell it, your reader, who by now will anyway have gotten that there were 10 and not 30 people involved in the study, will instead think either (i) that you *didn't know* that a higher user involvement would have been better to back your claims or (ii) that you intentionally want to *hide* information or even *cheat*. None of these are good for you, and for sure worse than telling straightaway. Keep this in mind.

Here some typical limitations of research. Check if any of them apply to your work:

- ☐ Small *sample size* (e.g., the number of users in the study or the amount of data collected in an experiment).

- ☐ For experiments that involve multiple *independent variables*, likely you will not have tested them all (e.g., in a crowdsourcing experiment, you fixed a reward for all experiments and did not study if that too affected your results).
- ☐ You may have *promised* something in the beginning of the thesis; if you didn't achieve everything either you drop the very promise or you mention it here as a limitation.
- ☐ When you collected data, there may have been some *bias* in the data (e.g., if you implement a prototype and do a user study yourself where there participants know that you actually implemented the software, they will give you biased answers, typically better ones).
- ☐ Collected data many have turned out being *incomplete* or of *lower quality* then initially expected. How does this impact your findings?
- ☐ Your prototype may have *crashed* or *not worked properly* in some experiments; it's important you tell the reader and explain possible implications of this on the validity of your conclusions.
- ☐ Due to time restrictions, you may have *not been able to complete* all experiments planned initially; again, explain the possible implications.
- ☐ People participating in a user study may have *dropped out* of the study, for whatever reason; if the reason is related to what you did or not did, you should mention it.
- ☐ Sometimes it is *not possible to compare* an own algorithm with other, similar algorithms, e.g., because their code is not available; this too may limit the viability of the findings.
- ☐ ...

8.4 Future Work

Finally, here you tell the reader which aspects you think would deserve further study or development. A good starting point for this is of course the list of limitations you just discussed. Not all of them may be worth investing more effort, but some will. The idea of this section is to identify where possible new effort should be invested, in order to make the work complete. Again, be frank and don't be afraid of identifying also new research directions. It's not you who will be doing what you propose here. It's meant for

the reader, the community. Everybody understands that after your defense you won't be working any longer on this project. It's all about suggesting future work, not telling that *you* will be doing it.

Bibliography

- [1] LIENS MSR INRIA)-Simon Lacoste-Julien (INRIA Paris Rocquencourt LIENS MSR INRIA) Aaron Defazio, Francis Bach (INRIA Paris Rocquencourt. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. 2014.
- [2] R. Harkreader C. Yang and G. Gu. Empirical evaluation and new design for fighting evolving twitter spammers. 2013.
- [3] Emilio Ferrara Alessandro Flammini Filippo Menczer Clayton A. Davis, Onur Varol. Botornot: A system to evaluate social bots, 2016.
- [4] Kyumin Lee, Brian David Eoff, and James Caverlee. Seven months with the devils: a long-term study of content polluters on twitter. 2011.
- [5] Francis Bach Mark Schmidt, Nicolas Le Roux. Minimizing finite sums with the stochastic average gradient. 2017.
- [6] Clayton A. Davis Filippo Menczer Alessandro Flammini Onur Varol, Emilio Ferrara. Online human-bot interactions: Detection, estimation, and characterization. 2017.
- [7] Cho-Jui Hsieh Xiang-Rui Wang Rong-En Fa, Kai-Wei Chang and Chih-Jen Lin. Liblinear: A library for large linear classification. 2017.
- [8] John S. and James L. Knight Foundation. Disinformation, 'fake news' and influence campaigns on twitter. 2018.
- [9] Marinella Petrocchi-Angelo Spognardi Maurizio Tesconi Stefano Cresci, Roberto Di Pietro. The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race. 2017.

Appendix A

User Manual

If you implemented a piece of software that is meant to be used by somebody else than you, then here you can provide a brief user manual that tells the target user how to use it. Part of this is the possible installation of the software and its operation and trouble shooting.

Appendix B

Dataset

If your work was based on a dataset that can be considered an output of the project, here you can describe it in detail.