School of Information Studies
Syracuse University

# IST565 DATA MINING
# Final Project Report

Instructor: Bei Yu

Student: David Forteguerre
SUID: 673608207
Email: dfortegu@syr.edu

Academic Year 2017 – 2018

**I. INTRODUCTION**

For this final project, I will explore the different levels of formality and registers of the French language using emails as the main form of communication.

English and Romance languages work very differently when it comes to addressing someone. In English, the pronoun *you* is the only option that can be used, whereas in Romance languages, there is a clear distinction in pronoun usage as someone may be addressed either informally (i.e., with the pronoun *tu* in most Romance languages) or formally (e.g. with the pronoun *vous* in French, *usted* in Spanish, *Lei* in Italian, etc.). Thus, the choice of pronoun in Romance languages is already a good indicator of the degree of formality of a given situation, even though it is not the only one. In France, a very high degree of formality is expected in many day-to-day situations where other countries and cultures would typically be more informal. For instance, in Spain, it is common to address a professor using *tu* instead of *usted*, to be on a first-name basis, and to greet each other saying *Hola*! (Hi!). In France, this type of behavior is unacceptable and would be seen as highly disrespectful. French professors may never be addressed with *tu* or by their first name; *vous* is the norm. Greetings also need to be more formal, so *Bonjour*! (Good morning!) would be used instead of *Salut*! (Hi!). The discrepancy in the level of formality required in different European countries for the same kind of situation can be explained by the fact that French culture has always shown a strong tradition of high respect and politeness, which has always been reflected in the language. This is certainly linked to the history of the language and also to the heavy influence of the French Academy, an institution that heavily regulates and protects the French language.

As a result, French is certainly the best language to choose for this analysis, culturally and linguistically speaking.

The purpose of this project is to discover the main linguistic patterns in the levels of formality of French emails. This is a typical classification task and several algorithms will be used in order to solve it. The classifiers will learn those patterns and ultimately output the main indicators of what makes an email either formal or informal.

Even though several data validation techniques will be used to evaluate the accuracy of the results (e.g. cross validation), I will also use my knowledge of French as a native speaker to evaluate the findings. The goal of this project is more educational than exploratory. Indeed, I plan to demonstrate a linguistic and sociolinguistic phenomenon through the use of algorithms to ultimately be able to explain it with real scientific evidence to a non-native speaker. From the beginning of the project, I already had a clear idea of what the algorithms were going to output, and my findings were eventually confirmed.
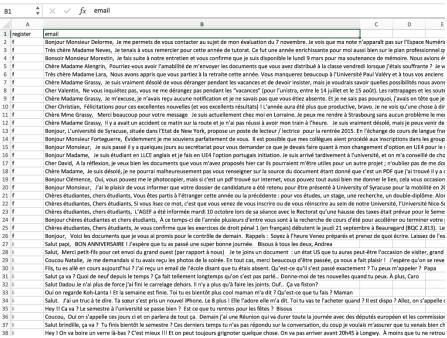
**II. DATA COLLECTION**

The data collection was the most tedious and time-consuming part of this project, but I wanted to make sure that I would have high-quality emails. When doing language analysis, there are key considerations to always have. The data must always reflect how the language is

actually used by a variety of people. Thus, I decided to ask my friends, their friends, as well as my family members to send me email samples. I made sure to ask people from all across France (mostly Southern France, Northern France, and Eastern France), and to ask people from different age ranges. I also included some of the emails I sent to my French 201 students at Syracuse University last year.

I collected **50** long and complete formal emails (more specifically university-related emails) and **50** informal emails (daily communications between family and friends.) As you know, the more data, the more accurate the algorithms output. Although it is often necessary to have a large amount of data for this type of task, I found the results to already be significant and very conclusive after running my analysis on 50 emails and then 100.

Even though a language often shows several registers and levels of formality (e.g. very formal, formal, everyday French, informal, slang), I decided to focus on the two major ones and create a data frame categorizing my data into binary bins: f (*formal*) vs. i (*informal*).

Below is a screenshot with fewer emails that demonstrates the organization of my data in Microsoft Excel.



emailsfrdata.csv[1]

Considering how difficult it can be to process data in Weka and following Professor Yu's recommendation, I decided to manually create the .arff file I would later import into Weka instead of attempting to convert my .csv file to an .arff file automatically.
However, in order to create that .arff file, I had to perform some preliminary steps on the data:
- I had to remove last names at the end of each email to preserve confidentiality.

---

[1] From just looking at the data, we can already notice some patterns at the beginning of the emails.

- I had to remove line breaks from the data to make sure all emails where in a paragraph form without line breaks.
- I also had to remove all apostrophes from the data. Apostrophes are very common in French as many grammatical words tend to be abbreviated when they come before a word that starts with a vowel. Oftentimes, those abbreviations are a grammatical requirement and not a register choice, unlike in English.

Abbreviations: **English vs. French**

| | 1st Element + | 2nd Element | Abbreviation | Required? |
|---|---|---|---|---|
| **French** | je (*I*) | ai (*have*) | j'ai (*I have*) | Required, as *je ai* would be grammatically incorrect |
| **English** | I | am | I'm | Not required. *I am* can also be used as a more formal option |

Only then was I able to create the .arff file. Below is a preview of the document.



emailsfrdata.arff

## III. DATA PREPROCESSING *(in Weka)*

Even though the dataset was now clean and in the right format, the data still needed to be preprocessed. Thus, I imported the emailsfrdata.arff file into Weka, and then selected the StringToWordVector filter (Unsupervised>attribute>StringToWordVector) in the "preprocess" tab. In the filter settings, I tuned some key parameters:
- I set **attributeIndices** to 2 (→ to apply the filter to the email variable").
- I set **lowerCaseTokens** to "true" (→ to make all capital letters lower case).

- I made sure that **minTermFrequencies** was set to "1" (default) (→ to remove any word that had only been used only *once* in all the emails combined).
- I set **normalizeDocLength** to normalize all data (→ to be able to see the word counts by clicking on each word after applying the filter).
- I set **stopwordsHandler** to my own stopwords file which I created (see below*). The file is called *stopwordsfr1.txt*.
- I set **tokenizer** to WordDelimiter (--> a basic tokenizer for this data).
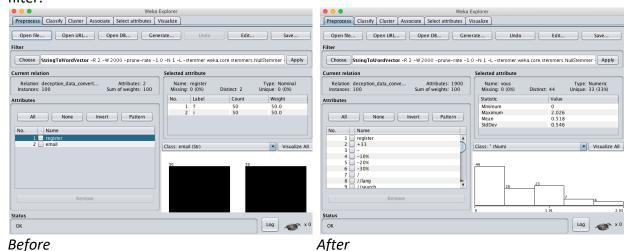- I set **wordsToKeep** to 2000 (instead of 1000, the default) (→ to have more data).

Finally, I applied the filter to the email variable whose content was processed and converted into word tokens.

*Stop words are words that are filtered out before or after processing of natural language data (text). Usually, these words are the most common non-content-bearing words in a language. However, there are no set stop words lists, as the list could vary depending on the nature of the analysis. In my case, I used R to get a list of French stop words. Below is my code:

```
# To get a list of French stopwords
library(tm)
FrenchStopWords <- stopwords("french")
print(FrenchStopWords, quote=FALSE)
```

I exported this list as a .txt file, made sure that each word was on a separate line, and finally went through all the words to make sure I would not be deleting any crucial words from the analysis. After going through the list, I found that "vous" (and its derivatives) and "tu" (and its derivatives) were in the list. Thus, I put a # sign before those words to make sure they would be kept in the dataset.

Below is a screenshot *before* applying the filter, and the output generated *after* applying the filter.



*Before*                                                    *After*

## IV. DATA ANALYSIS *(in Weka)*

In data mining, a very common technique is to run several algorithms on the same dataset and tune their parameters in order to see which one works best and outputs the highest accuracy. Indeed, each algorithm has unique properties and tends to be best suited for a specific set of tasks only.

For this classification task, I decided to run three algorithms: (1) Naïve Bayes, (2) Support Vector Machines, and (3) Decision Trees in order to see which would work best.

### 1. Naïve Bayes[2]

I first ran the Naïve Bayes Multinomial algorithm on the dataset.[3] The data was ready to be analyzed, thus I selected the NaiveBayesMultinomial algorithm in the "classify tab", and made sure that the target variable was set to "register" (i.e. what we are trying to predict.) I used 10-fold cross validation in order to evaluate the model.[4]

The results are as follows:

| Naïve Bayes algorithm | |
|---|---|
| **Default settings:** | **Correctly Classified Instances      100         100   %** |
| | Observation: all instances have been classified accurately. |

Let's now look at the confusion matrix for the result. It shows that all instances have been classified accurately.

```
=== Confusion Matrix ===

 a  b   <-- classified as
50  0 |  a = f
 0 50 |  b = i
```

This perfect accuracy may point out that there are some obvious patterns in the data that make it very easy for the algorithm to find the answer. By taking a look at the probabilities of given words, and keeping in mind how pronouns work in French, we can observe major differences. Let's take the example of "tu" (informal *you*).

| "tu" | | |
|---|---|---|
| | **f** | **i** |
| Probability output | 3.982777131170316E-4 | 0.01294848936653945 |
| Value rounded-up | **0.0003** | **0.012** |

---

[2] Note that several Naïve Bayes algorithms exist. NB Multinomial is the one that is particularly designed for text categorization.

[3] There are two NB Multinomial algorithms in Weka. One is called NaiveBayesMultinomial and needs the data to be preprocessed and filtered as input, and the other is called NaiveBayesMultinomialText—it actually takes care of the preprocessing and filtering section. I tried both algorithms. For NaiveBayesMultinomialText, I had to click "Undo" in the preprocess tab to undo the output of my StringToWordVector, and then tune the parameters within the algorithm in a similar way before running it. **NaiveBayesMultinomial** actually showed a much higher accuracy. Thus, NaiveBayesMultinomialText was disregarded.

[4] Cross-validation is a technique to evaluate predictive models by partitioning the original sample into a training set to train the model, and a test set to evaluate it. In k-fold cross-validation, the original sample is randomly partitioned into k equal size subsamples.

As you can see, the informal pronoun "tu" is much more likely to appear in an informal email than in a formal one. So far, our hypothesis has been verified. Let's go even further and use a different algorithm.

### 2. SVMs

The second step was to run Support Vector machines on the dataset. Thus, I selected the SMO algorithm (under "functions") and made sure that the target variable was still set to "register." Just like for NB, I used 10-fold cross validation in order to evaluate the model.
The results are as follows:

| SVM algorithm | |
|---|---|
| **Default settings:** | Default kernel (= PolyKernel)<br>**Correctly Classified Instances      100        100    %**<br>Observation: all instances have been classified accurately. |

Once again, all emails have been classified accurately. The confusion matrix was the same as the NB algorithm matrix. Using another algorithm and getting similar results really emphasizes that there must be something in the data that betrays the level of formality of each email. It was not even necessary to tune the algorithm parameters for this specific task to try and get a higher accuracy, as the outputs were already 100% accurate (which is unusual.)
Once again, let's take a look at the attribute weights to see if we can discover any patterns correlated with the usage of pronouns in French.

| Pronoun | Weight |
|---|---|
| "vous" | +      -0.3534 * (normalized) vous |
| "tu" | +       0.3058 * (normalized) tu |
| Note that a positive weight predicts an informal email in this output, whereas a negative weight predicts a formal email. | |

### 3. Decision Trees

Finally, I decided to run the Decision Tree algorithm. My thought process was that this type of algorithm would be very well suited for this type of task where only a few key elements seem to play a determining role in our predictions. Indeed, it is possible to easily visualize the patterns found by this algorithm thanks to an actual tree output.
To build this model, I selected Weka's J48 algorithm and made sure the target variable was still set to "register". Once again, I used 10-fold cross validation in order to evaluate the model.
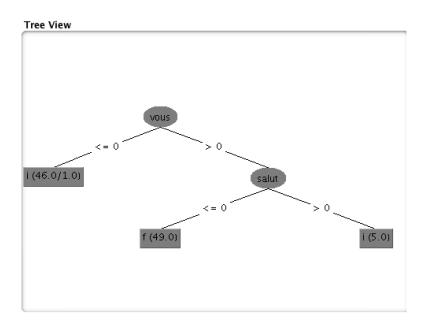The results are as follows:

| J48 algorithm | |
|---|---|
| **Default settings:** | **Correctly Classified Instances      99        99    %**<br>Observation: once again, the accuracy obtained is exceptionally high. |

Below is the confusion matrix of the output. We can see that a formal email has been mistakenly classified as informal.

```
=== Confusion Matrix ===

  a   b   <-- classified as
 49   1 |   a = f
  0  50 |   b = i
```

Here is the visualization of the tree generated by the model:

**Tree View**



This tree confirms the hypothesis introduced at the beginning of this project. "Vous" (formal *you*) seems to be *by far* the linguistic element that betrays the level of formality of French emails.

This tree suggests that if there is no "vous" in a given email, it is likely to be informal.[5] On the other hand, if there is one or more "vous", then the email could very well be formal or not. That is when the analysis becomes even more interesting and informative. The second key indicator in the data I collected seems to be the informal greeting "salut" (*hi*). If an email contains one or more "vous" <u>and</u> contains one or more "salut", then it is likely to be informal. Indeed, this greeting is typically used in very informal situations and would never be used in a formal email (such as a university-related communication). However, if an email contains one or more "vous" and <u>does not</u> contain any "salut", then it is very likely to be formal.
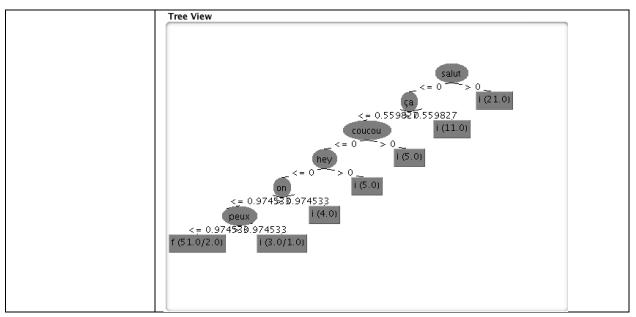
**Conclusion**

---

[5] Note that the error found must come from the fact that "vous" is not only used to mean *you* (formal, to address 1 person), but also *you* (no specific register, to address several people) (= "you guys").

The algorithms used were all helpful in this analysis. As suggested, pronoun usage seems to be key in French to determine the register of a given communication.

However, I decided to push this analysis further and take out "tu" and "vous" from the data to see how the algorithms would perform and discover what the next patterns would be. In order to do so, I created a second stop words list (*stopwordsfr2.txt*) and removed the # signs I had placed before "vous" (and its derivatives) and "tu" (and its derivatives). I then reimported the dataset into Weka, and reapplied the StringToWordVector filter, this time selecting the brand new stop words list. I double-checked to make sure that the pronouns "vous" and "tu" had been removed from the data, and then ran the three algorithms again, using 10-fold cross validation.

Below is a quick summary of the results.

| 1. NB | **Correctly Classified Instances** | **98** | **98** | **%** |
|---|---|---|---|---|
| | **Incorrectly Classified Instances** | **2** | **2** | **%** |
| | `=== Confusion Matrix ===`<br><br>`  a  b   <-- classified as`<br>`48  2 |   a = f`<br>` 0 50 |   b = i` | | | |
| 2. SVMs | **Correctly Classified Instances** | **100** | **100** | **%** |
| | `=== Confusion Matrix ===`<br><br>`  a  b   <-- classified as`<br>`50  0 |   a = f`<br>` 0 50 |   b = i` | | | |
| 3. Decision Trees | **Correctly Classified Instances** | **83** | **83** | **%** |
| | **Incorrectly Classified Instances** | **17** | **17** | **%** |
| | `=== Confusion Matrix ===`<br><br>`  a  b   <-- classified as`<br>`43  7 |   a = f`<br>`10 40 |   b = i` | | | |

**Tree View**



Running the algorithms again without *tu* and *vous* tells us a lot about the data and the performance of the algorithms to solve this type of task. Indeed, it is now becoming obvious that SVMs seem to be best suited for the task, NB coming in second position, and decision trees coming last. However, we also notice that the accuracy given by decision trees is 83%. Even though it is now much lower, it remains a very high accuracy overall.

This extra step showed us that some other key words also play a role in determining register in French, and these are:
1. **Salut** (also found in the previous results), which means *hi* and is typically used in informal contexts.
2. **Ça**, which means *this* or *that*. However, *ça* is the informal abbreviation of *cela*. Thus, this demonstrative is used in informal contexts as well and serves as the second best indicator.
3. **Coucou**, which means *hi*. It is another informal greeting.
4. **Hey**, which also means *hi/hey*. Again, it is an informal greeting.
5. **On**, which is the informal "nous" (i.e. the informal *we*). This shows that another pronoun (other than "vous" and "tu") is also a very good indicator.
6. […] Then, we start getting a list of inflected verbs and conjugations (such as "peux"), which is not relevant to this analysis.

As a final step, I decided to use GainRatio on the data (including "tu" and "vous" again) and list the top-20 features. In order to do so, I selected the InfoGainAttributeEval evaluator in the "select attribute" tab. Below is the list of the rankings:

0.8601   1829 **vous**
0.4142   1739 **tu**
0.2508   1538 **salut**

```
0.2461    1081 ne
0.2266    1824 votre
0.1953     980 madame
0.1853     272 bonjour
0.1821    1050 monsieur
0.1821    1853 ça
0.1692    1139 on
0.1441     450 cordialement
0.1441      77 a
0.1441    1433 remercie
0.1198    1823 vos
0.108     1264 plus
0.108      266 bisous
0.108     1673 te
0.108     1771 va
0.0964     704 faire
0.0964     451 coucou
```

This list confirms our findings. Pronouns and greetings played a key role in this analysis, and those features helped our classifiers learn different patterns in the French email data.


**V. CONCLUSIONS**

To conclude, we have seen that our algorithms have successfully carried out this text classification task.

- First of all, we have seen that the difference of pronoun usage between "tu" and "vous" plays a key role in determining the level of formality of French emails. However, those pronouns were not the only factors, and greetings came next. Indeed, some greetings such as "salut" (*hi*!) are only found in very informal contexts, and also play a role in betraying the register of a given email.

- This analysis also emphasizes a very important difference between French (and more generally, Romance Languages) and English. Romance languages tend to have a wider variety of pronouns, and each pronoun usually conveys very specific ideas and information. Indeed, there is no distinction between *you* (singular) and *you* (plural) in English, which is why this type of analysis would be irrelevant to the English language.

- Finally, this analysis not only highlights two linguistic patterns, but also emphasizes a major difference between American and French cultures. If French culture was as informal and relaxed as American culture (where students often address their professors with an informal tone and often write to them using informal greetings such as "hi" or "hello"), this analysis would not have been as conclusive. Thus, we also learned that the education system in France always requires a high level of formality.

**Real life application**

       As we learned that French culture is much stricter than American culture when it comes to formality, one could imagine a situation where a French university wants to create and implement its own focused inbox that would filter out any email that is not work-related. Even though this concept already exists (Microsoft Outlook does offer the option to have a focused inbox to only keep important emails), it is unlikely that the algorithms developed by Microsoft only look at formality levels in order to classify important and trivial emails due to the informal nature of American culture. Indeed, an important email (e.g., a professor's reminder) may still have an informal tone in the U.S. This would not be the case in France, and it would be much easier to implement this type of tool using formality as the main (and only?) factor.

**APPENDIX 1**

**EXAMPLES OF CONJUGATIONS IN FRENCH SHOWING PRONOUN USAGE**

### 1. Regular verbs

|  | PARLER<br>*to speak* | FINIR<br>*to finish* |
|---|---|---|
| **je** | parle | finis |
| **tu** | parles | finis |
| **il/elle/on** | parle | finit |
| **nous** | parlons | finissons |
| **vous** | parlez | finissez |
| **ils/elles** | parlent | finissent |

### 2. Irregular verbs

|  | ÊTRE<br>*to be* | AVOIR<br>*to have* |
|---|---|---|
| **je** | suis | (j') ai |
| **tu** | es | as |
| **il/elle/on** | est | a |
| **nous** | sommes | avons |
| **vous** | êtes | avez |
| **ils/elles** | sont | ont |
|  | *p. 29* | *p. 72* |

**APPENDIX 2**

**TU OR VOUS USAGE IN FRENCH?**

Below you will find guidelines for non-native speakers that teach you when to use *tu* or *vous* in French. Even though this image is somewhat satirical and comes from *BuzzFeed*, it is probably the best and most accurate visualization I have ever seen as a native speaker.