

Damegender Manual: Counting Males and Females in Internet Communities

for version 0.4, 30 Jul 2022

David Arroyo Menéndez (davidam@gmail.com)

This manual is for Damegender (version 0.4, 30 Jul 2022).

Copyright © 2020 David Arroyo Menéndez

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.3 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the Appendix A entitled "License".

Table of Contents

1	Introduction	1
1.1	Reproducible Research	1
2	Installation	3
3	Commands	4
4	Statistics	10
4.1	Measuring success and error	10
4.2	Principal Component Analysis (PCA)	14
4.2.1	Counting features in names	14
4.2.2	Choosing components	15
5	Use Cases	18
5.1	Introduction	18
5.2	Counting males and females in Debian	18
5.3	Counting males and females in Linux Kernel	20
5.4	Counting males and females in Forbes	20
5.5	Deciding for males and females in images	22
5.6	Webscraping and Damegender (counting scholars)	22
5.7	Counting males and females in a git repository	25
5.8	Counting males and females in Maps	26
5.9	Gender gap in science	27
6	Secondary Sources about the Gender Gap	28
6.1	Gender Inequality in the World	28
6.2	Gender Inequality in STEM	30
6.3	Gender Inequality in Free Software	31
7	Theoretical Frameworks	33
7.1	Philosophies about software, market, freedom and gender	33
7.2	Multiculturalism, Interculturalism	36
7.3	Feminism, Ecofeminism and Intersectionality	38
7.4	Gender	39
8	Conclusions	40
Acknowledgments		41

Further reading	42
Appendix A License.....	45
Appendix B Photos	53
Index.....	57

1 Introduction

Damegender is a gender detection tool from the name coded by David Arroyo MEnéndez (DAME). See “*Damegender: Writing and Comparing Gender Detection Tools*”, [Further reading], page 42.

The gender detection tools from the names are being used usually with commercial APIs. But many countries has been doing efforts in the last years contributing names and gender datasets with Open Data Licenses. So, this software is collecting this effort in an industrial way and giving new original solutions to classify gender from the name (we are using Machine Learning algorithms predicting names that is not appearing in our datasets).

Damegender is giving measures to compare in any moment our solution with the commercial APIs. So, the user can understand when it's useful to invest money or not depending of the dataset. Damegender allows to the users to download a big number of names from a csv file.

This software is written oriented to tests. So, you can check the right behaviour of the software with python tests for the classes and methods and with shell tests for the python commands.

Damegender is using Perceval helping to count males and females in a lot of Internet Communities (wikis, mailing lists, software repositories, bug tracking systems, ...). “*Perceval: software project data at your will*”, [Further reading], page 42. This manual to show source for count males and females in different contexts (Ex: `git2gender.py`, `mail2gender.py`).

This software is taking into account the power to predict nations and ethnicity from the surnames (Ex: `surname.py`, `surnameincountries.py` and `ethnicity.py`).

This book starts explaining the installation. Later, how to install Damegender. After, it shows Damegender from the commands. Next, it explains the statistical maths concepts to use Damegender with good results. The use cases showed allows to imagine a lot of applications about Damegender. We are giving some data sources about gender gap. And finally, the theoretical frameworks are being showed to put data into discourses. We are added conclusions as summary and further reading for extend the interest.

1.1 Reproducible Research

Damegender has been written thinking on make easy the reproducible research. It's possible using:

- Free Software for Damegender and the Damegender dependencies.
- Free documentation for Damegender (this manual, FAQs, and a website)
- Open Data with sources written.
- Good practices about programming (POO, using standards, unity tests and user tests)

The are some interesting goal towards reproducible research. That's more link between code and data in a way where you can download all data and execute all code. These tasks has not been reached due to:

- To retrieve data from Internet becomes slow and you can suffer cuts in the Internet connection.

- To train machine learning models and sometimes execute machine learning models with a good number of names is slow. Then it's not practical for build a PDF make all work of code and data linked.

So in the spectrum of reproducibility (Peng, 2011) the damegender papers and damegender manual will be in the point of partially linked and executable code and data.

2 Installation

Possible Debian/Ubuntu dependencies:

```
$ sudo apt-get install python3-nose-exclude python3-dev dict
dict-freedict-eng-spa dict-freedict-spa-eng dictd
```

Now, to install damegender with python package:

```
$ python3 -m venv /tmp/d
$ cd /tmp/d
$ source bin/activate
$ pip install --upgrade pip
$ pip3 install damegender
$ cd lib/python3.5/site-packages/damegender
$ python3 main.py David
```

To install API extra dependencies:

```
$ pip3 install damegender[apis]
```

To install mailing lists and repositories extra dependencies:

```
$ pip3 install damegender[mails_and_repositories]
```

To install all possible dependencies

```
$ pip3 install damegender[all]
```

Currently you can need an API key from:

- <https://store.genderize.io/documentation>
- <https://gender-api.com>
- <https://www.nameapi.org/>
- <https://v2.namsor.com/NamSorAPIv2/sign-in.html>

To configure your api key you can execute:

```
$ python3 apikeyadd.py
```

3 Commands

You must start to check tests to understand that all is OK:

```
$ cd src/damegender
$ ./testsbycommands.sh          # It must run for you
$ ./testsbycommandsextralocal.sh # You will need all dependencies
# with: $ pip3 install damegender[all]
$ ./testsbycommandsperceval.sh   # Only about perceval
$ ./testsbycommandsextraapis.sh # You will need all dependencies
$ ./testmergeinterfiles.sh      # Only about merge datasetes
$ ./testsorig2.sh               # Checking orig2.py in several countries■
```

You can continue checking python tests:

Execute all tests:

```
$ nosetests3 tests
```

Another way to execute all tests:

```
$ ./run-unit-test.sh
```

Execute one file:

```
$ nosetests3 tests/test_basics.py
```

Execute one test:

```
$ nosetests3 tests/test_basics.py:TestBasics.test_indexing
```

If you are in a fresh installation, perhaps you want to regenerate by your own risk some files generated by me, such as, ML models to understand how it has been generated:

```
$ python3 postinstall.py
```

Another task related is merging datasets with names and frequency, for instance, you can merge 10 females of Denmark with 10 females of Germany with:

```
python3 mergeinterfiles.py
--file1="files/names/names_inter/dkfemales10.csv"
--file2="files/names/names_inter/defemales10.csv"
--output=files/tests/testdkde-$(date "+%Y-%m-%d-%H").csv
```

You can download and regenerate all csv files by your own risk with:

```
$ ./testsorig2.sh           # Checking orig2.py in several countries■
$ ./regenerate-inter-files.sh # Regenerate inter names files
$ ./regenerate-intersurnames.sh # Regenerate inter surnames files
$ ./regenerate-malefemale-files.sh # Regenerate male and female files
```

We have results applying ML to our dataset cached in json files. You can regenerate these json files with:

```
$ ./regenerate-ml-json.sh
```

You can find an big list of commands to execute this shell scripts. Now a detailed execution of some selected examples:

The first command to learn is `main.py`. You can play now with this command:

```
# Detect gender from a name (INE is the dataset used by default)
$ python3 main.py David
```

```

David gender is male
363559 males for David from INE.es
0 females for David from INE.es

# Detect gender from a name only using machine learning
$ python3 main.py Agua --ml=nltk
Agua gender is female
0 males for Agua from INE.es
0 females for Agua from INE.es

# Detect gender from a name (all census and machine learning)
$ python3 main.py David --verbose
365196 males for David from INE.es
0 females for David from INE.es
1193 males for David from Uruguay census
5 females for David from Uruguay census
26645 males for David from United Kingdom census
0 females for David from United Kingdom census
3552580 males for David from United States of America census
12826 females for David from United States of America census
David gender predicted with nltk is male
David gender predicted with sgd is male
David gender predicted with svc is male
David gender predicted with gaussianNB is male
David gender predicted with multinomialNB is male
David gender predicted with bernoulliNB is male
David gender predicted with forest is male
David gender predicted with tree is male
David gender predicted with mlp is male

```

The first Free Software for gender detection tool was created in C language program and you can look for a python version with the name genderguesser. Some people was working in a Free dataset called name_dict.txt with 48500 names. I want to give thanks to this effort with `nameincountries.py` due to the good work organizing many names in different countries.

```

$ python3 nameincountries.py David
grep -i " David " files/names/nam_dict.txt > files/grep.tmp
males: ['Albania', 'Armenia', 'Austria', 'Azerbaijan', 'Belgium',
'Bosnia and Herzegovina', 'Czech Republic', 'Denmark', 'East Frisia',
'France', 'Georgia', 'Germany', 'Great Britain', 'Iceland', 'Ireland',
'Israel', 'Italy', 'Kazakhstan/Uzbekistan', 'Luxembourg', 'Malta',
'Norway', 'Portugal', 'Romania', 'Slovenia', 'Spain', 'Sweden',
'Swiss', 'The Netherlands', 'USA', 'Ukraine']
females: []
both: []

```

To count gender from a git repository:

```
$ python3 git2gender.py
```

```
https://github.com/chaoss/grimoirelab-perceval.git
--directory="/tmp/clonedir"
The number of males sending commits is 15
The number of females sending commits is 7
```

You can see a verbose output using the Spanish dataset (`--language=es`) for males and females with:

```
$ python3 git2gender.py https://git.drupalcode.org/project/orgmode.git
--directory=/tmp/orgmode --show=all --verbose --language=es
You are not using ml the process is not very slow, but perhaps
you are not finding good results
The number of males sending commits is 2
The list of males sending commits is:
['David Arroyo Menendez', 'David Arroyo']
David Arroyo Menéndez <davidam@es.gnu.org> (67 commits)
David Arroyo Menendez <davidam9@riseup.net> (49 commits)
David Arroyo Menéndez <davidam@gmail.com> (20 commits)
David Arroyo Menendez <david.arroyo@edoctores.com> (10 commits)
David Arroyo Menendez <davidam@es.gnu.org> (14 commits)
David Arroyo7 <davidam@es.gnu.org> (13 commits)
David Arroyo7 <davidam@gnu.org> (10 commits)
The number of females sending commits is 1
The list of females sending commits is:
['Miriam']
Miriam <miriam@xxxxxx.es> (23 commits)
The number of people with unknown gender sending commits is 0
The list of people with unknown gender sending commits is:
[]
```

To count gender from a mailing list:

```
$ cd files/mbox
$ wget -c
http://mail-archives.apache.org/mod_mbox/httpd-announce/201706.mbox
$ cd ../..
$ python3 mail2gender.py
http://mail-archives.apache.org/mod_mbox/httpd-announce/
You are not using ml the process is not very slow, but perhaps you are
not finding good results
```

```
The number of males sending mails is 24
The number of females sending mails is 2
The number of people with unknown gender sending mails is 5
```

You can execute a verbose output with:

```
$ python3 mail2gender.py
http://mail-archives.apache.org/mod_mbox/httpd-announce/
--verbose --show=all
You are not using ml the process is not very slow, but perhaps you are not
```

```

finding good results
The number of males sending mails is 24
The list of males sending mails is:
['Jim <jim@xxxxxx.es>', 'Jacob <jchampion@xxxxx.org>',
'DENNIS <balaranpillai@xxxxx.com>', '"Leonard (Jira)" <jira@xxxxx.org>',
'"Roy" <jira@xxxxx.org>', '"Roman (Jira)" <jira@xxxxx.org>',
'"Bertrand" <jira@xxxxx.org>', '"Mark (Jira)" <jira@xxxxx.org>',
'"Justin (Jira)" <jira@xxxxx.org>', '"Simon (Jira)" <jira@xxxxx.org>',
'"Chris (Jira)" <jira@xxxxx.org>', 'Jan <lahoda@xxxxx.com>',
'"Michael (Jira)" <jira@xxxxx.org>', '"Ralph (Jira)" <jira@xxxxx.org>',
'"Jens" <jensg@xxxxx.org>', 'Mark <markt@xxxxx.org>',
'"Ryan (Jira)" <jira@xxxxx.org>', 'Ismaël (Jira) <jira@xxxxx.org>',
'"Shane (Jira)" <jira@xxxxx.org>', '"Kevin A. (Jira)" <jira@xxxxx.org>',
'"Gordon (Jira)" <jira@xxxxx.org>', 'Gary <garydgregory@xxxxx.com>',
'"Owen" <owen.omalley@xxxxx.com>', '"Sheng (Jira)" <jira@xxxxx.org>']

The number of females sending mails is 2
The list of females sending mails is:
['Riya <hellen.serviceweb@xxxxxx.com>',
'"Hannah (Jira)" <jira@xxxxx.org>']

The number of people with unknown gender sending mails is 5
The list of people with unknown gender sending mails is
[ '"SimpaticoTech" <web.info@xxxxxx.it>',
'Simpatico <web.info@xxxxxxxxxx.it>',
'gmcdonald@xxxxx.org',
'Hen <bayard@xxxxx.org>',
'"Jean (Jira)" <jira@xxxxx.org>']

```

With Damegender is possible to guess the authors in an article about a newspaper with:

```
$ python3 newspaper2gender.py https://elpais.com/espana/catalunya/2021-11-23/el-presup
```

Perhaps you don't know a name, but you have obtained an free key for an API to retrieve it:

```
$ python3 api2gender.py Leticia --surname="Martin" --api=namsor
female
scale: 0.99
```

If you want to know the gender of a good number of names you can download results from an api and save in a file with downloadjson.py.

```
$ python3 downloadjson.py --csv=files/names/min.csv --api=genderize
$ cat files/names/genderizefiles_names_min.csv.json
```

Now we are going to learn some commands for measure the successful of our solution:

```
$ python3 accuracy.py --csv=files/names/min.csv
#####
NLTK!!
Gender list: [1, 1, 1, 1, 2, 1, 0, 0]
Guess list: [1, 1, 1, 1, 0, 1, 0, 0]
Dame Gender accuracy: 0.875
$ python3 confusion.py --csv="files/names/partial.csv" --api=nameapi
--jsongDownloaded="files/names/nameapifiles_names_partial.csv.json"
```

A confusion matrix C is such that $C_{i,j}$ is equal to the number of observations known to be in group i but predicted to be in group j. If the classifier is nice, the diagonal is high because there are true positives Nameapi confusion matrix:

```
[[ 3, 0, 0]
 [ 0, 15, 1]]
$ python3 errors.py --csv="files/names/all.csv" --api="genderguesser"
Gender Guesser with files/names/all.csv has:
+ The error code: 0.22564457518601835
+ The error code without na: 0.026539047204698716
+ The na coded: 0.20453365634192766
+ The error gender bias: 0.0026103980857080703
```

You can generate a lot of logs about errors, accuracies and/or confusion:

```
$ ./logs-accuracies.sh
$ ./logs-confusion.sh
$ ./logs-errors.sh
$ ./logs-count.sh
```

Perhaps you are interested on reproduce experiments to determine features:

```
$ python3 infofeatures.py
# To determine number of components
$ python3 pca-components.py --csv="files/features_list.csv"
# To understand the weight between variables for a target
$ python3 pca-features.py
```

Now we can go to play with surnames:

```
$ python3 surname.py Gil --total=es
There are 140004 people using Gil in Spain

$ python3 surname.py Lenin --total=us
There are 837 people using Lenin in United States of America
```

```
$ python3 ethnicity.py Smith
In United States of America the percentages about the race
of Smith surname is:
White: 73.35
Black: 22.22
Hispanic: 1.56
Asian Pacific Indian American: 0.40
American Indian and Alaska Native: 0.85
Various races: 1.63
```

You can download sources chosen by damegender and generate the damegender csv files from any country with:

```
$ python3 orig2.py es --download
```

To generate json files, such as, a rest service is downloaded you can use the next python command

```
$ python3 csv2jsonapirest.py files/names/names_inter/dkfemales10.csv --outdir="files/t
```

4 Statistics

In the last chapter, we were learning to execute some commands such as `accuracy.py`, `confusion.py`, or `errors.py`, but perhaps you need to understand more theory about statistics to understand why this commands is being interesting for you.

4.1 Measuring success and error

To guess the sex, we have an true idea (ex: female) and we obtain a result with a method (ex: using an api, querying a dataset or with a machine learning model). The guessed result could be male, female or perhaps unknown. To remember some definitions about results about this matter:

True positive is to find a value guessed as true if the value in the data source is positive.

True negative is to find a value guessed as true if the the value in the data source is negative.

False positive is to find a value guessed as false if the the value in the data source is positive.

False negative is to find a value guessed as false if the the value in the data source is negative.

So, we can find a vocabulary for measure true, false, success and errors. We can make a summary in the gender name context about mathematical concepts:

Precision is about true positives divided by true positives plus false positives

$$\frac{\text{femalefemale} + \text{malemale}}{\text{femalefemale} + \text{malemale} + \text{femalemale}}$$

Recall is about true positives divided by true positives plus false negatives.

$$\frac{\text{femalefemale} + \text{malemale}}{\text{femalefemale} + \text{malemale} + \text{malefemale} + \text{femaleundefined} + \text{maleundefined}}$$

Accuracy is about true positives divided by all.

$$\frac{\text{femalefemale} + \text{malemale}}{\text{femalefemale} + \text{malemale} + \text{malefemale} + \text{femalemale} + \text{femaleundefined} + \text{maleundefined}}$$

The **F1 score** is the harmonic mean of precision and recall taking both metrics into account in the following equation:

$$2 * \left(\frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \right)$$

In Damengender, we are using `accuracy.py` to apply these concepts. Take a look to the execution:

```
$ python3 accuracy.py --api="damegender" --measure="f1score"
--csv="files/names/partialnoudefined.csv"
--jsondownloaded=files/names/partialnoudefined.csv.json
#####
# Damegender!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
```

```

        1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list: [1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1,
             1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0]
Damegender f1score: 0.9090909090909091

$ python3 accuracy.py --api="damegender" --measure="recall"
--csv="files/names/partialnoundefined.csv"
--jsondownloaded=files/names/partialnoundefined.csv.json
##### Damegender!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
              1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list: [1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1,
             1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1]
Damegender recall: 1.0

$ python3 accuracy.py --api="damegender" --measure="accuracy"
--csv="files/names/partialnoundefined.csv"
--jsondownloaded=files/names/partialnoundefined.csv.json
##### Damegender!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
              1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list: [1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1,
             1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1]
Damegender accuracy: 0.8571428571428571

$ python3 accuracy.py --api="genderguesser" --measure="accuracy"
--csv="files/names/partialnoundefined.csv"
--jsondownloaded=files/names/partialnoundefined.csv.json
##### Genderguesser!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
              1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list: [1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1,
             1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1]
Genderguesser accuracy: 0.8571428571428571

$ python3 accuracy.py --api="genderguesser" --measure="precision"
--csv="files/names/partialnoundefined.csv"
--jsondownloaded=files/names/partialnoundefined.csv.json
##### Genderguesser!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
              1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list: [1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1,
             1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1]
Genderguesser precision: 0.9090909090909091

$ python3 accuracy.py --api="genderguesser" --measure="recall"
--csv="files/names/partialnoundefined.csv"

```

```
--jsondownloaded=files/names/partialnoundefined.csv.json
#####
Genderguesser!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
              1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list:  [1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1,
              1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0]
Genderguesser recall: 1.0

$ python3 accuracy.py --api="genderguesser" --measure="f1score"
--csv="files/names/partialnoundefined.csv"
--jsondownloaded=files/names/partialnoundefined.csv.json
#####
Genderguesser!!
Gender list: [1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
              1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1]
Guess list:  [1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1,
              1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1]
Genderguesser f1score: 0.9090909090909091
```

Error coded is about the true is different than the guessed:

```
(femalemale + malefemale + maleundefined + femaleundefined) /
(malemale + femalemale + malefemale +
+ femalefemale + maleundefined + femaleundefined)
```

Error coded without na is about the true is different than the guessed, but without undefined results.

```
(maleundefined + femaleundefined) /
(malemale + femalemale + malefemale +
+ femalefemale + maleundefined + femaleundefined)
```

Error gender bias is to understand if the error is bigger guessing males than females or vice versa.

Weighted error is about the true is different than the guessed, but giving a weight to the guessed as undefined.

```
(femalemale + malefemale +
+ w * (maleundefined + femaleundefined)) /
(malemale + femalemale + malefemale + femalefemale +
+ w * (maleundefined + femaleundefined))
```

In Damegender, we have coded `errors.py` to implement several definitions in different API.

The confusion matrix creates a matrix about the true and the guess. If you have this confusion matrix:

```
[[ 2, 0, 0]
 [ 0, 5, 0]]
```

It means, I have 2 females true and I've guessed 2 females and I've 5 males true and I've guessed 5 males. I don't have errors in my classifier.

```
[[ 2   1   0]
 [ 2 14   0]]
```

It means, I have 2 females true and I've guessed 2 females and I've 14 males true and I've guessed 14 males. 1 female was considered male, 2 males was considered female.

In Damegender, we have coded `confusion.py` to implement this concept:

```
python3 confusion.py --csv=files/names/min.csv
--api=damegender --jsondownloaded=files/names/min.csv.json
A confusion matrix C is such that Ci,j is equal to the number of
observations known to be in group i but predicted to be in group j.
If the classifier is nice, the diagonal is high because there are
true positives
Damegender confusion matrix:

      M   F   U
M  [[ 5,  0,  0 ]
F  [ 0,  1,  0 ]]
```

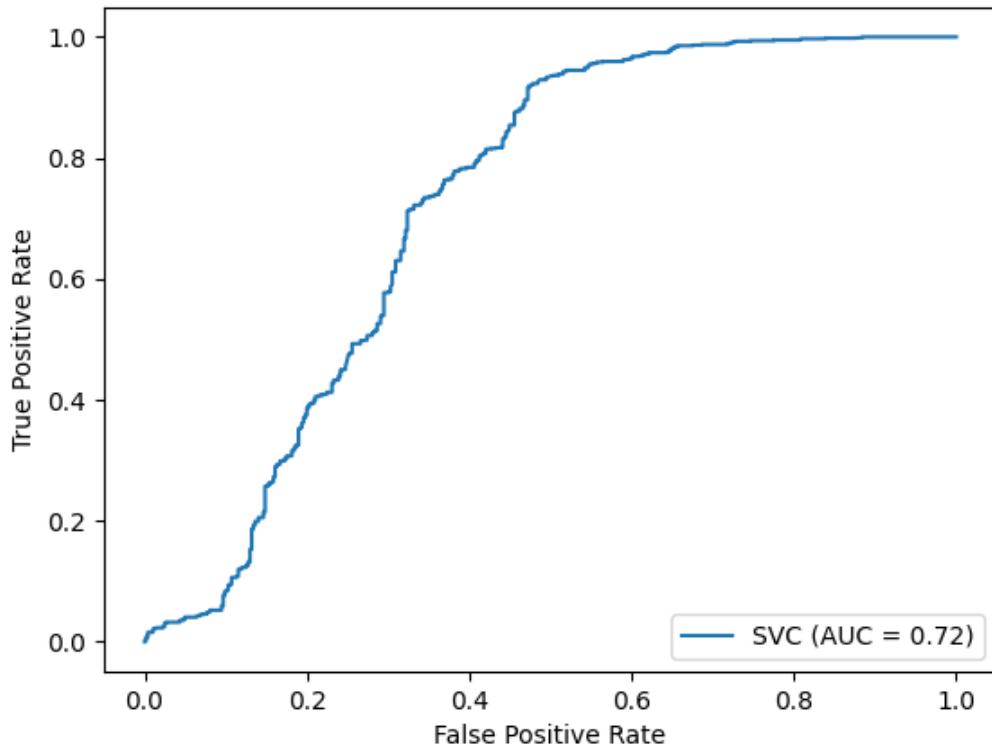
Remember that we can retrieve the json file from several apis having the names to be guessed in the csv file with `downloadjson.py`:

```
$ python3 downloadjson.py --csv="files/names/min.csv" --api="genderapi"
```

Similar to confusion is ROC (Receiver Operating Characteristic) is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. The ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings.

In Damegender, you can use ROC relative to machine learning algorithms with the next command:

```
$ python3 roc.py svc
```



4.2 Principal Component Analysis (PCA)

4.2.1 Counting features in names

We have developed a script `infofeatures.py` with our datasets to visualize data about some features chosen by us.

```
$ python3 infofeatures.py ine
```

Take a look to the results with the different datasets:

Dataset	Letter A	Last Letter A	Last Letter O	Last Letter Consonant	Last Letter Vocal	First Letter Consonant	First Letter Vocal
Uruguay (females)	0.816	0.456	0.007	0.287	0.712	0.823	0.177
Uruguay (males)	0.643	0.249	0.062	0.766	0.234	0.771	0.228
Australia (females)	0.922	0.588	0.033	0.272	0.728	0.772	0.228
Australia (males)	0.818	0.03	0.269	0.57	0.43	0.763	0.237
Canada (females)	0.659	0.189	0.005	0.591	0.408	0.838	0.161

Canada (males)	0.752	0.22	0.025	0.54	0.456	0.818	0.181
Spain (females)	0.922	0.588	0.03	0.271	0.728	0.772	0.228
Spain (males)	0.818	0.03	0.268	0.569	0.43	0.763	0.236
United Kingdom (females)	0.825	0.374	0.013	0.322	0.674	0.765	0.235
United Kingdom (males)	0.716	0.036	0.039	0.78	0.218	0.799	0.2
USA (females)	0.816	0.456	0.007	0.287	0.712	0.823	0.177
USA (males)	0.643	0.02	0.061	0.765	0.234	0.84	0.159

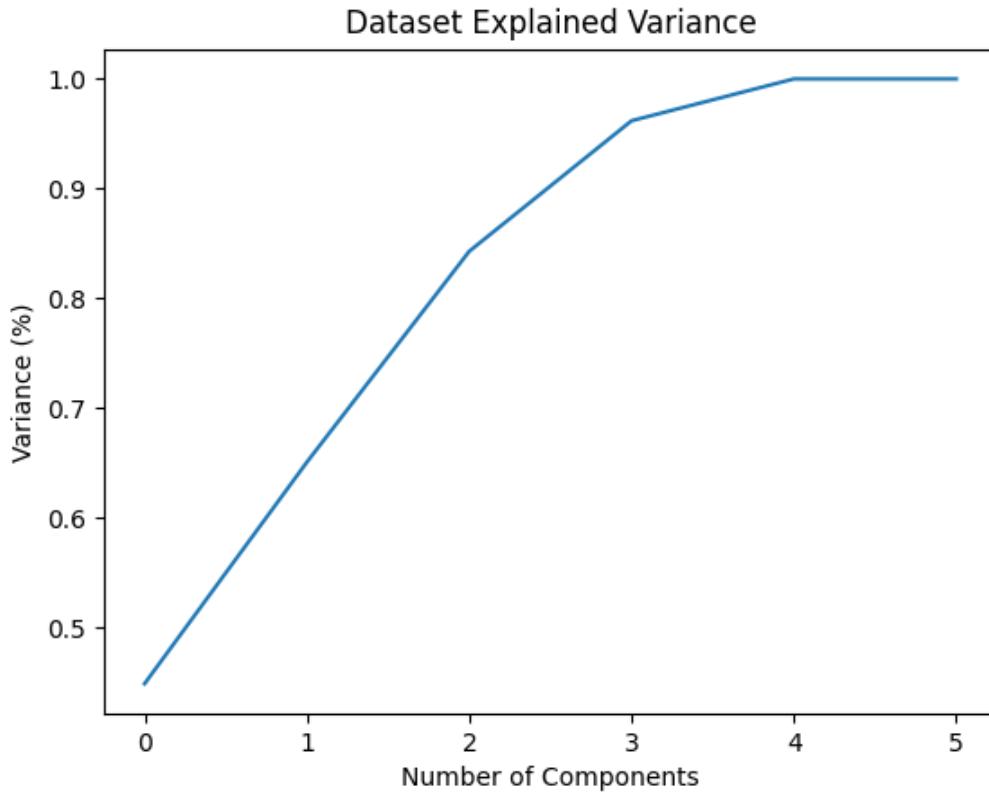
The countries where the main language is Spanish (Uruguay + Spain) and English (USA + United Kingdom + Australia) are having very similar variation with the features chosen between males and females with these datasets (to remember is the datasets extracted from official statistics provided by the states). Canada, a country where french is the main language has different rules with this features.

The letter a is varying 0.2 from males to females in (USA and Uruguay) and 0.1 from males to females (United Kingdom, Australia and Spain). The last letter a is varying 0.5 from males to females in (Australia, Spain) around 0.4 in (USA, United Kingdom) and 0.2 in Uruguay. The last letter o from females to males is varying 0.2 in (Spain, Australia) and is equal in (Uruguay, USA, United Kingdom). For the last letter consonant all countries is giving the result that is for males, with results from 0.2 to 0.5: Uruguay and USA (0.5), United Kingdom (0.4), Australia and Spain (0.3). So last letter vocal is reverse the last letter consonant. First letter consonant or first letter vocal is a non significant feature due to so similar results in English and Spanish.

Surely, the rules it's a coincidence but we think that is a coincidence between languages due to that there are a good number of names to think different.

4.2.2 Choosing components

After, to choose features for our machine learning task, we can understand if this features makes sense with Principal Component Analysis. We have written 2 scripts for this task `pca-components.py` and `pca-features.py`. With `pca-components.py` we are giving a csv (files/features_list.csv, files/features_list_no_cat.csv, ...) and the output is an image where we can visualize a curve to determine when this curve stops the growth the number of components.



In the image, we can see that the curve stops the growth in the fourth component.

When you know the components you can execute `pca-features.py` so:

```
$ python3 pca-features.py --categorical=both --components=4
The json file is created in files/pca.json
The html file is created in files/pca.html
```

first_letter	last_letter	last_letter_a	first_letter_vocal	last_letter_vocal	last_letter_consonant	target component
-0.2080025204	-0.3208958517	0.2352509625	0.2113242731	0.6095269139	-0.6095269139	-0.1035071139
-0.6037951881	0.5174873789	-0.4252467151	0.4278794455	0.0388287435	-0.0388287435	-0.0265942125
0.1049343046	0.1158117877	-0.2867605971	-0.3473950734	0.0901034539	-0.0901034539	-0.8697264971
0.2026467275	0.3142402839	0.630802294	0.5325769702	-0.1291229841	0.1291229841	-0.3811720011

To simplify and to learn, we can observe this analysis without letters. In this analysis, we can observe 4 components.

The first component is about if the last letter is vocal or consonant. If the last letter is vocal we can find a female and if the last letter is a consonant we can find a male.

The second component is about the first letter. The last letter is determining females and the first letter is determining males.

The third component is not giving relevant information.

The fourth component is giving the last_letter_a and the first_letter_vocal is for females.

5 Use Cases

5.1 Introduction

There are many research studies count males and females in specific communities such as Twitter, Stack Overflow, ... A specific community has some clues to determine male or female, for example, in Twitter you can observe the photo, nickname, real name, ...

In this chapter we are going to apply the concepts to determine gender in real situations observing where the gender is provided.

5.2 Counting males and females in Debian

In the Debian community all members must have a gpg key to collaborate, so we can count males and females from the keyring. With gpg commands you can import a the debian keyring and dump the debian keyring in a csv file.

```
$ rsync -az --progress keyring.debian.org::keyrings/keyrings/ .
```

We have generated a script to count males and females:

```
~/git/damegender/src/damegender$ python3 count-debian-gender.py
Perhaps you need wait some minutes. You can take a tea or coffee now
Debian males: 795
Debian females: 24
```

In the dump of the Debian keyring dataset we have divided name, surname and email in different fields. So, it's easy detect the name, although some names has several emails.

We have chosen the United States of America dataset and we are using the method name_frec to decide for male or female in the row. Take a look to the source:

```
import csv
import unicodedata
import unidecode
from pprint import pprint
import re
from app.dame_gender import Gender
from app.dame_utils import DameUtils

du = DameUtils()
g = Gender()

result=""
dm = []

with open('files/debian-maintainers-gpg-2020-04-01.csv') as csvfile:
    reader = csv.reader(csvfile, delimiter=',', quotechar='|')
    aux = ""
    cnt = 0
    for row in reader:
        cnt = cnt +1
```

```

        if (aux != row[0]):
            dm.append(row[0])
        aux = row[0]

print("Perhaps you need wait some minutes. You can take a tea or coffe now")■

females = 0
males = 0
for rowdm in dm:
    if (int(g.name_frec(str(rowdm.upper()), 'us')['females']) > int(g.name_frec(str(rowdm.upper()), 'us')['males'])):
        females = females + 1
    else:
        males = males + 1

print("debian males: %s" % males)
print("debian females: %s" % females)

 csvfile.close()

```

The advantage using the method `name_frec` is about to understand how you are deciding male or female in the script counting males and females. In this script the decision is simple: a name is male if there are more males than females and female if there are more females than males.

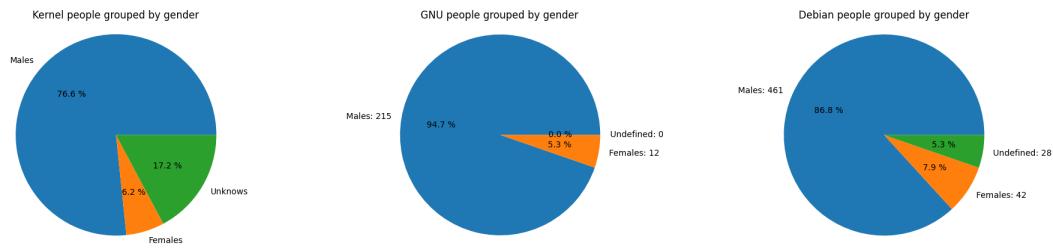
The United States of America dataset is a good choice for Free Software communities, due to that this communities is based on English as main language and United States of America is a leader country in software development. United States of America hosts people from different countries due to migrations towards good companies and universities.

In general, you can choose `csv2gender.py` to count males, females and unknowns in a csv file. For example, doing this:

```
$ python3 csv2gender.py --first_name_position=0
files/debian-maintainers-gpg-2020-04-01.csv --verbose
```

But, a research must understand the source, too.

In the next diagram we can see (78.4% of males, 5.4% of females and 16.2% of unknowns).



We can retrieve the names of unknowns and to decide about the name in the sense of retrieve the gender from a commercial API (genderapi, genderize, namsor, nameapi, ...) or

to classify the name as a software company or a bug. For now, the Open datasets contributed by states about names are very good but not all countries has the idea of contribute the names. Remember that you can download names from a csv file with `downloadjson.py`:

```
$ python3 downloadjson.py --api=genderize --csv=files/names/min.csv
```

5.3 Counting males and females in Linux Kernel

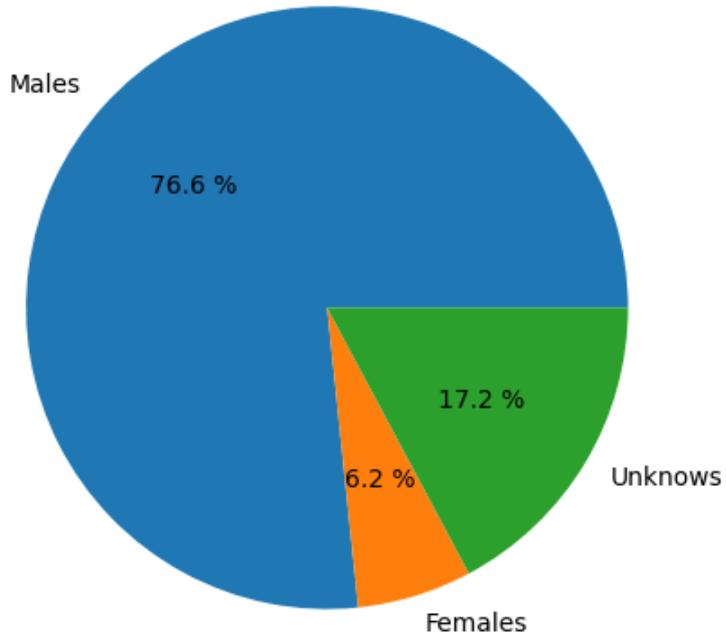
When I'm writing this book the Linux Kernels maintainers appears in <https://www.kernel.org/doc/html/latest/process/maintainers.html>. Then I have downloaded the file and applied a single command:

```
cat maintainers.html | w3m -dump -T text/html  
| grep "Mail:" > maintainers.txt
```

You can makes fixes to this command from GNU/Emacs or with shell scripting. Later, you can apply:

```
$ python3 count-kernel.py
```

Kernel people grouped by gender



5.4 Counting males and females in Forbes

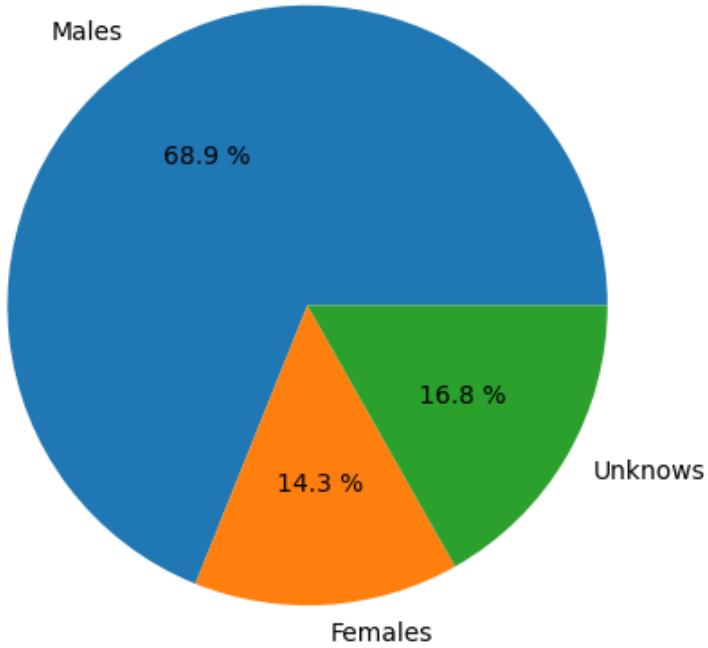
In the second example, we are using guess without machine learning instead of `name_frec`. If you are using `guess` you are trusting on `damegender` to take the decision, but perhaps you are not agree.

Please take a look about our guess method in the current state:

```
def guess(self, name, binary=False, *args, **kwargs):
    # guess list method
    dataset = kwargs.get('dataset', 'es')
    # guess method to check names dictionary
    guess = ''
    name = unidecode.unidecode(name).title()
    name.replace(name, "")
    dicc = self.name_frec(name, dataset)
    m = int(dicc['males'])
    f = int(dicc['females'])
    if ((m == 0) and (f == 0)):
        if binary:
            guess = 2
        else:
            guess = "unknown"
    elif (m > f):
        if binary:
            guess = 1
        else:
            guess = "male"
    elif (f > m):
        if binary:
            guess = 0
        else:
            guess = "female"
    else:
        if binary:
            guess = 2
        else:
            guess = "unknown"
    return guess
```

We are using the Spanish dataset by default and the rest is the same idea that in the last script: more people using the name.

Top 119 Forbes people grouped by gender



5.5 Deciding for males and females in images

There are many free software tools for decide gender in images files. We can use these tools to decide gender about images from Twitter, GitHub, ... We have selected the next tool:

```
$ git clone https://github.com/davidam/damefaces
$ cd damefaces/bin
$ python3 damefaces.py girl1.jpg
```

5.6 Webscraping and Damegender (counting scholars)

Sometimes, we can reach the database of names from a website, for example, we can retrieve a list of academics from Spain thanks to webometrics and the next script:

```
from lxml import html
import requests

print("Introduce an url from webometrics, for example,
      https://www.webometrics.info/en/GoogleScholar/Spain")

import argparse

parser = argparse.ArgumentParser()
parser.add_argument("url", help="display the gender")
args = parser.parse_args()
```

```
page = requests.get(args.url)
tree = html.fromstring(page.content)

academics = tree.xpath('//tr/td/a/strong/text()')

print('Academics: %s' % academics)
```

If you have retrieved the list of names in a file `files/scientifics.txt`, you could count males and females with the next script called `count-scientifics.py`:

```
import csv
import unicodedata
import unidecode
import re

from pprint import pprint
from app.dame_gender import Gender
from app.dame_utils import DameUtils
from ast import literal_eval
from app.dame_sexmachine import DameSexmachine

du = DameUtils()
g = Gender()
s = DameSexmachine()

with open('files/scientifics.txt') as f:
    mainlist = [list(literal_eval(line)) for line in f]

l = mainlist[0]

ll = []
for i in l:
    ll.append(i.split())

ten = ll[0:10]
hundred = ll[0:100]
thousand = ll[0:1000]

x = 0
y = 0
males = 0
females = 0
for j in hundred:
    if (len(j[0]) == 1):
        x = x + 1
    else:
```

```
sex = g.guess(j[0], binary=False)
y = y +1
if (sex == "male"):
    males = males + 1
elif (sex == "female"):
    females = females + 1

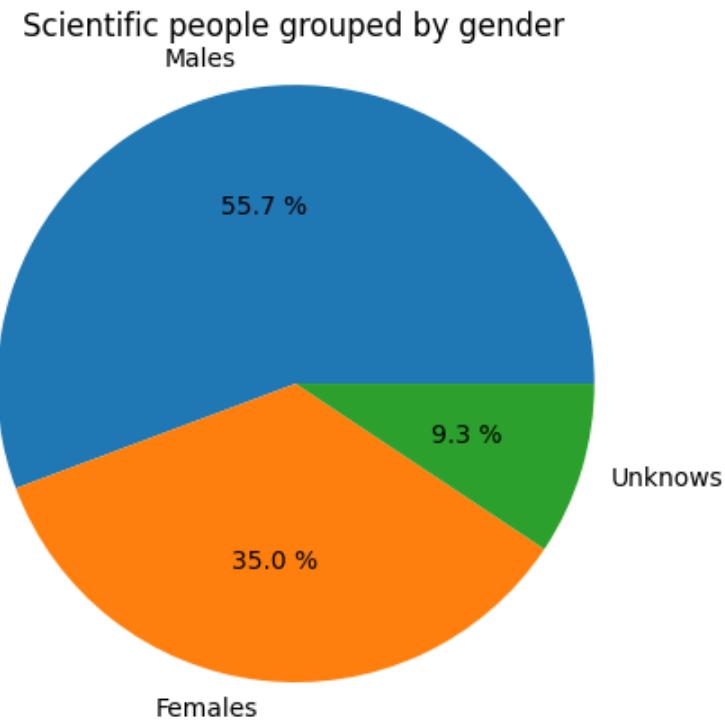
print("Number of scientifics with a single letter as first name: %s" % x)
print("Number of scientifics with the first name normal: %s" % y)
print("Number of females scientifics: %s" % females)
print("Number of males scientifics: %s" % males)

for j in thousand:
    if (len(j[0]) == 1):
        x = x + 1
    else:
        sex = g.guess(j[0], binary=False)
        y = y +1
        if (sex == "male"):
            males = males + 1
        else:
            females = females + 1

print("Number of females scientifics: %s" % females)
print("Number of males scientifics: %s" % males)
```

And the results are:

```
Number of females scientifics: 31425
Number of males scientifics: 47945
```



So, the percentage of academic people classified as females in Spain is bigger than Free Software people classified as females. Having a gender gap in both situations.

5.7 Counting males and females in a git repository

We can think a simple version of `git2gender.py`:

```
from app.dame_sexmachine import DameSexmachine
from app.dame_perceval import DamePerceval
from app.dame_utils import DameUtils
import sys
import argparse
parser = argparse.ArgumentParser()
parser.add_argument("url", help="Uniform Resource Link")
parser.add_argument('--directory')
parser.add_argument('--version', action='version', version='0.1')
args = parser.parse_args()
if (len(sys.argv) > 1):
    ds = DameSexmachine()
    du = DameUtils()
    dp = DamePerceval()
    l1 = dp.list_committers(args.url, args.directory)
    l2 = du.delete_duplicated(l1)
    l3 = du.clean_list(l2)
```

```
females = 0
males = 0
unknowns = 0
for g in ls:
    sm = ds.guess(g, binary=True)
    if (sm == 0):
        females = females + 1
    elif (sm == 1):
        males = males + 1
    else:
        unknowns = unknowns + 1

print("The number of males sending commits is %s" % males)
print("The number of females sending commits is %s" % females)
```

Try to execute this script:

```
$ python3 git2gender.py https://github.com/davidam/davidam.git
--directory="/tmp/clonedir"
The number of males sending commits is 3
The number of females sending commits is 0
```

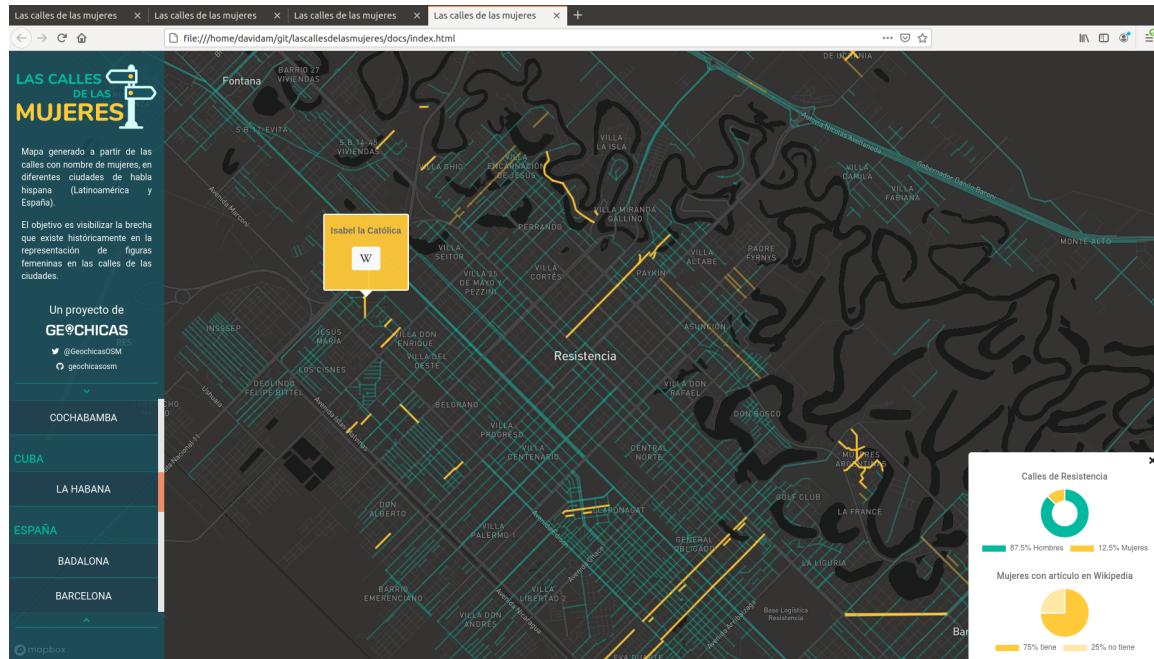
5.8 Counting males and females in Maps

Las calles de las mujeres is a project in NodeJS to display web streets with names about females using MapBox. You can download the project with:

```
$ git clone https://github.com/geochicasosm/lascallesdelasmujeres
```

It's licensed with a Creative Commons License (CC-BY-SA)

<https://creativecommons.org/licenses/by-sa/4.0/>.



It's giving statisticals about how many men and women has Wikipedia article, too.

It has a strong community created by females (GeoChicas) <https://geochicas.github.io/> related with OpenStreetMap.

5.9 Gender gap in science

A good visual job can be found in <https://lukeholman.github.io/genderGap/>. This work is based on R retrieves the data from genderize and arxiv. The gap is especially large in authorship positions associated with seniority, and prestigious journals have fewer women authors. Additionally, they estimates that men are invited by journals to submit papers at approximately double the rate of women. Wealthy countries, notably Japan, Germany, and Switzerland, had fewer women authors than poorer ones. It concludes that the STEMM gender gap will not close without further reforms in education, mentoring, and academic publishing. There are a paper with a full explanation: “*The gender gap in science: How long until women are equally represented?*”, [Further reading], page 42.

6 Secondary Sources about the Gender Gap

When a social researcher starts a new work, the first step is set an objective about the project with subobjectives, that's define the problem to solve with a methodology (quantitative, qualitative, or mixed). The second step is about sample decision who is the people, the population going to give us the data. The third step is about selection strategies about retrieve data, analysis and to show results. (Source: *Técnicas Cualitativas de Investigación Social*) (Source: *La encuesta una perspectiva general metodológica*) results. See “*Técnicas Cualitativas de Investigación Social*”, [Further reading], page 42,

To read secondary sources about the objective and subobjectives of the project is to read previous works of another people (papers, books, data, news, ...) in science we refer to the state of the art or similar.

You can read secondary sources in different steps of a social research work, although is a task specially suggested in the first steps.

The objective of this chapter is to show some secondary sources about the gender gap in general, in STEM and in the Free Software. The idea is to reduce time to the people using Damegender understanding this kind of things. After of this, you can compare gender gap with the local problem: Kernel, Gnu/Linux distribution, Stack Overflow, Twitter, Forbes, Science, etc.

6.1 Gender Inequality in the World

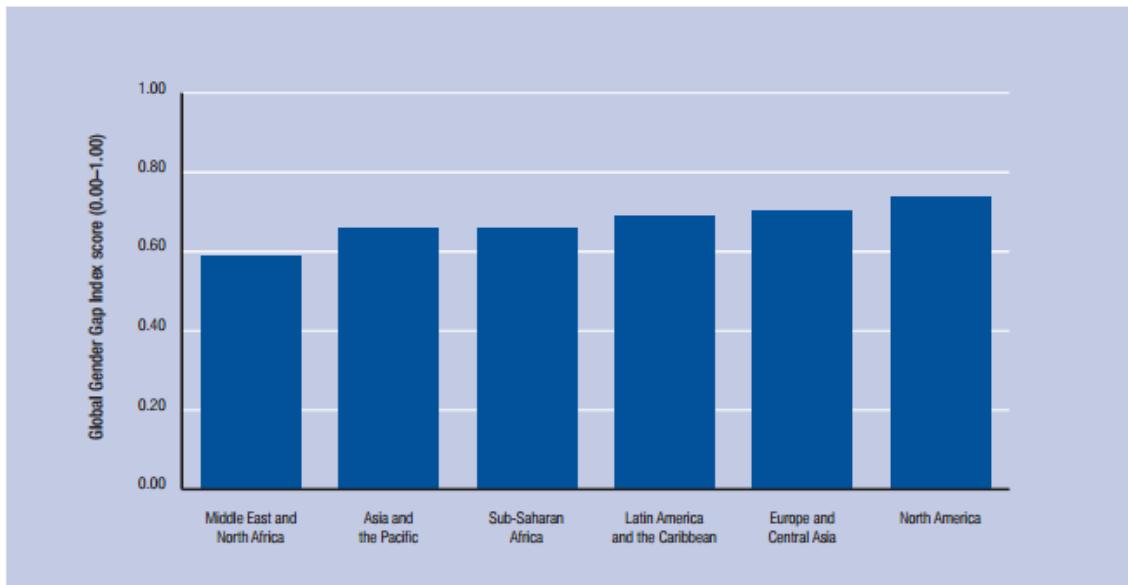
Gender gap or gender inequality is the idea that men and women are not equal and that gender affects an individual's living experience. These differences arise from distinctions in biology, psychology, and cultural norms. Some of these types of distinctions are empirically grounded while others appear to be socially constructed. Studies show the different lived experience of genders across many domains including education, life expectancy, personality, interests, family life, careers, and political affiliations. Gender inequality is experienced differently across different cultures. (Source: Wikipedia, 2020)

The women is underrepresented in the labour world (among adults aged from 25 to 54 has stagnated over the past 20 years, standing at 31 percentage points. The gender pay gap exists, too, so the women are paid 16% less than men. Share of women and men with an account at a financial institution is 65% of the total in women and 72% of the total in men. 31% of young women aged 15 to 24 are not in education, employment or training in 2020, more than double rate for young men (14%). Violence against women is 18% of ever-partnered women aged 15 to 49 experienced sexual and/or physical violence by an intimate partner in the previous 12 months.¹

¹ <https://www.unwomen.org>



In the next image you can observe a graphic about gender gap in several continents (Source: Global Gender Gap Index 2012):

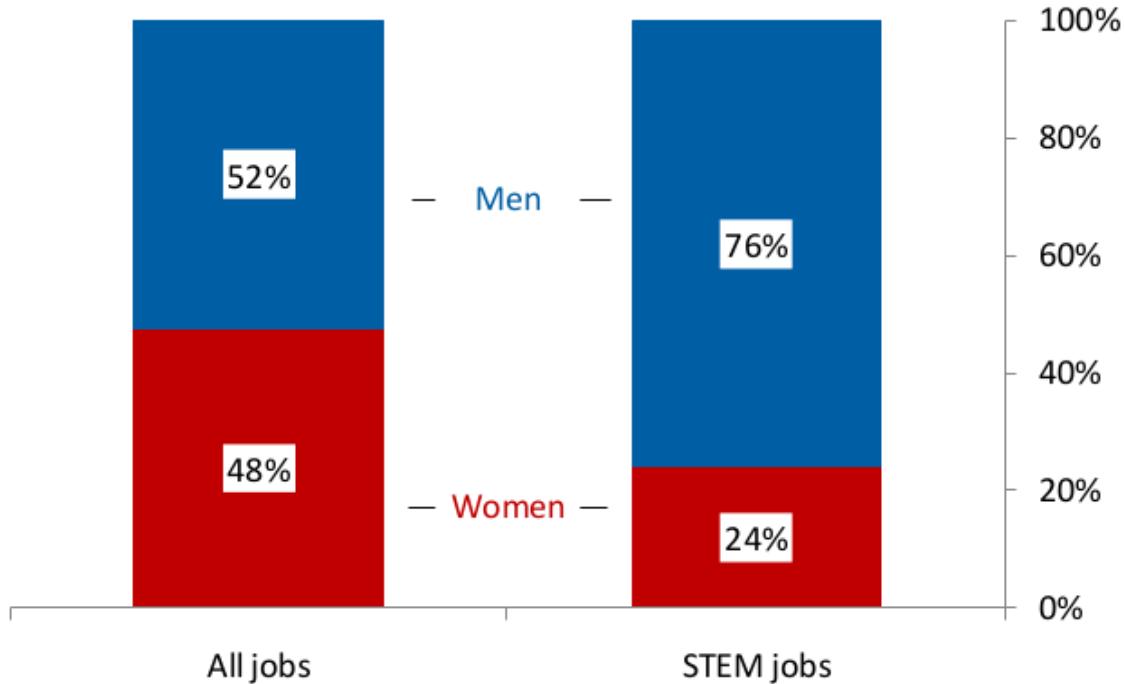
Figure 2: Regional performance on the Global Gender Gap Index 2012

Source: Global Gender Gap Index 2012; details of regional classifications in Appendix B.
Scores are weighted by population; population data from the World Bank's *World Development Indicators* (WDI) online database 2011, accessed July 2012.

From the best score to the worst score. We can find: North America with the best score, Europe and Central Asia, Latin America and the Caribbean, Sub-Saharan Africa, Asia and the Pacific, and the worst score Middle East and South Africa.

6.2 Gender Inequality in STEM

In the graphic we can understand gender gap in stem in 2009 and compare with gender gap in the market:

Figure 1. Gender Shares of Total and STEM Jobs, 2009

Source: ESA calculations from American Community Survey public-use microdata.

Note: Estimates are for employed persons age 16 and over.

So the gender gap in STEM is bigger than the labour market in general.

6.3 Gender Inequality in Free Software

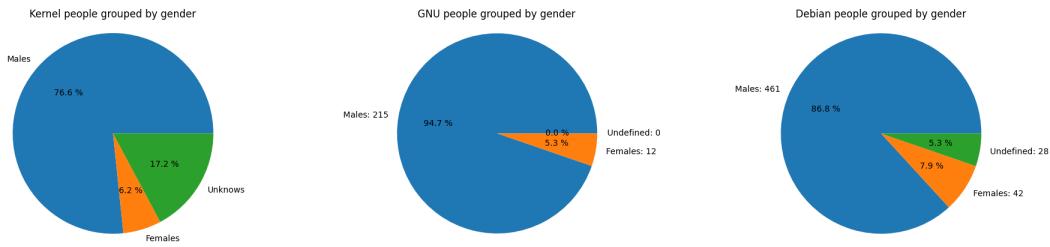
The gender gap in the Free Software world is so high we are presenting preliminar due to that we can reduce the percentage of unknowns to male or female. But we can observe that the gender gap is bigger in the Free Software world than in STEM.

The context of the operatings systems is not feminist by different reasons. In the Free Software world there are many males developing the software. But in another Operating Systems there are another problems about the male domination, for example, Microsoft will create the most rich men in the world for many years². Apple was classified as the most valuable company in the world³. By intersectionality, to create companies more powerful than states is bad for the democracy and in a world where exists the gender gap the change to the gender equity is a pressure of values.

A Linux operating system such as we find in a CD of Ubuntu, RedHat or Suse is composed by the Kernel (Linux) and GNU programs. Later, it's organized in a distribution, such as Ubuntu and Debian. We have choosen Debian as use case and we are counting males females in Linux, too.

² <https://www.forbes.com/profile/bill-gates/?sh=79472717689f>

³ <https://forbes.es/listas/32245/apple-la-marca-mas-valiosa-del-2017/>



In Debian, the number of males is 461 (86.8%), the number of females is 42 (7.9%) and the number of names undefined is 28.

In GNU, the number of males is 215 (94.7%) and the number of females is 12 (5.3%).

In the Kernel, the number of males is 1387 (72.5%) and the number of females is 62 (3.2%). There are a big number of undefined, so we must debug these data, but we have not time in this version of the manual.

7 Theoretical Frameworks

If you want do a social research for a quantitative study, such as, count males and females you can:

- To generate objectives about the research study
- To read and to understand previous works.
- To choose some theoretical framework.
- Perhaps, with a qualitative study (interviews, focus groups, ...) you could understand better the problems, the specific vocabulary used by the people, the reasons about the decisions, ...
- To retrieve data with Damegender or make surveys (online, offline, ...)
- An analysis quantitative with maths, graphics and interpretations
- Conclusions

A theoretical framework consists of concepts and, together with their definitions and reference to relevant scholarly literature, existing theory that is used for your particular study. In the last chapters you have learnt to count males and females, but you need give meaning to the words that you are using about your gender study. That is the point in this chapter.

We present some theoretical frameworks that you can use as example in your works:

- Philosophies about software market and freedoms
- Interculturalism and Multiculturalism
- Feminism, Ecofeminism and derivatives
- Gender terms and philosophies

7.1 Philosophies about software, market, freedom and gender

There are different philosophies developing software and we are counting males and females in Internet, so the floor is the software in this world. If we must analyze gender in a country the ideology is changing in the place where you are. In the software world is the same problem. So, we are giving the vocabulary and the philosophy for speak about software and ideologies.

The proprietary software is the most common idea for the common people, operating systems such as Microsoft Windows or Mac OS. If you are using software with proprietary licenses, the source files will be containing copyright notes such as:

```
# Copyright (C) 2020 David Arroyo Menéndez

# Author: David Arroyo Menéndez <davidam@gmail.com>
# Maintainer: David Arroyo Menéndez <davidam@gmail.com>

# All rights reserved
```

This idea is associated to big companies leading the market but any people can use this philosophy. The criticism appears with Richard Stallman about privacy and lack of freedom

to the academic people, or hackers (people who knows read and write software and they do it for his objectives or global objectives). I could to say the monopoly is too strong with this license and the current social inertia and now nobody can change the market, we need another licenses to preserve the free market with an ethical strategy for startups and students.

Richard Stallman defines the Free Software with four freedoms: (0) to run the program, (1) to study and change the program in source code form, (2) to redistribute exact copies, and (3) to distribute modified versions. See “Free Software Free Society”, [Further reading], page 42,

This idea to build software as a social good and motivated by ethical values. The solution is to apply the GPL license and to request to GNU to include the software.

The copyright note in GNU would be similar to:

```
; ; This software is free software: you can redistribute it and/or modify
; ; it under the terms of the GNU General Public License as published by
; ; the Free Software Foundation, either version 3 of the License, or
; ; (at your option) any later version.

; ; This software is distributed in the hope that it will be useful,
; ; but WITHOUT ANY WARRANTY; without even the implied warranty of
; ; MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
; ; GNU General Public License for more details.

; ; You should have received a copy of the GNU General Public License
; ; along with GNU Emacs. If not, see <https://www.gnu.org/licenses/>.
```

On opposition the Open Source movement believes in free licenses, but they thinks that the software is business and they want to develop Free Software by economy, so they prefer change the word Free Software by Open Source claiming their philosophy.

They redefines the Free Software Definition by the Open Source Definition¹.

1. Free Redistribution

The license shall not restrict any party from selling or giving away the software as a component of an aggregate software distribution containing programs from several different sources. The license shall not require a royalty or other fee for such sale.

2. Source Code

The program must include source code, and must allow distribution in source code as well as compiled form. Where some form of a product is not distributed with source code, there must be a well-publicized means of obtaining the source code for no more than a reasonable reproduction cost, preferably downloading via the Internet without charge. The source code must be the preferred form in which a

¹ <https://opensource.org/osd>

programmer would modify the program. Deliberately obfuscated source code is not allowed. Intermediate forms such as the output of a preprocessor or translator are not allowed.

3. Derived Works

The license must allow modifications and derived works, and must allow them to be distributed under the same terms as the license of the original software.

4. Integrity of The Author's Source Code

The license may restrict source-code from being distributed in modified form only if the license allows the distribution of "patch files" with the source code for the purpose of modifying the program at build time. The license must explicitly permit distribution of software built from modified source code. The license may require derived works to carry a different name or version number from the original software.

5. No Discrimination Against Persons or Groups

The license must not discriminate against any person or group of persons.

6. No Discrimination Against Fields of Endeavor

The license must not restrict anyone from making use of the program in a specific field of endeavor. For example, it may not restrict the program from being used in a business, or from being used for genetic research.

7. Distribution of License

The rights attached to the program must apply to all to whom the program is redistributed without the need for execution of an additional license by those parties.

8. License Must Not Be Specific to a Product

The rights attached to the program must not depend on the program's being part of a particular software distribution. If the program is extracted from that distribution and used or distributed within the terms of the program's license, all parties to whom the program is redistributed should have the same rights as those that are granted in conjunction with the original software distribution.

9. License Must Not Restrict Other Software

The license must not place restrictions on other software that is distributed along with the licensed software. For example, the license must not insist that all other programs distributed on the same medium must be open-source software.

10. License Must Be Technology-Neutral

No provision of the license may be predicated on any individual technology or style of interface.

In the point six, we find the conflict with the feminist theories due to the positive discrimination is a good idea to reach gender equity.

The GNU philosophy has the same problem explained on a different way. Only the Free Software is a good idea, if the software is not Free Software, then it's Proprietary Software (the bad idea to avoid).

In Damegender, we want to deliver Free Software by practical reasons released with GPLv3 through pypi.org and github.com the very popular sites to distribute Free Software written in Python. But we understand that we can make positive discrimination in the development in favor to the women as an experiment with this copyright note in the development branch:

```
# You can share, copy and modify this software if you are a woman or you
# are David Arroyo Menéndez and you include this note.
```

This book has been developed with this license, too.

7.2 Multiculturalism, Interculturalism

The term multiculturalism has a range of meanings within the contexts of sociology, of political philosophy, and of colloquial use. In sociology and in everyday usage, it is a synonym for "ethnic pluralism", with the two terms often used interchangeably, for example, a cultural pluralism in which various ethnic groups collaborate and enter into a dialogue with one another without having to sacrifice their particular identities. It can describe a mixed ethnic community area where multiple cultural traditions exist (such as New York City or Trieste) or a single country within which they do (such as Switzerland, Belgium or Russia). Groups associated with an indigenous, aboriginal or autochthonous ethnic group and settler-descended ethnic groups are often the focus.

In reference to sociology, multiculturalism is the end-state of either a natural or artificial process (for example: legally-controlled immigration) and occurs on either a large national scale or on a smaller scale within a nation's communities. On a smaller scale this can occur artificially when a jurisdiction is established or expanded by amalgamating areas with two or more different cultures (e.g. French Canada and English Canada). On a large scale, it can occur as a result of either legal or illegal migration to and from different jurisdictions around the world (for example, Anglo-Saxon settlement of Britain by Angles, Saxons and Jutes in the 5th century or the colonization of the Americas by Europeans, Africans and Asians since the 16th century).

In reference to political science, multiculturalism can be defined as a state's capacity to effectively and efficiently deal with cultural plurality within its sovereign borders. Multiculturalism as a political philosophy involves ideologies and policies which vary widely. It has been described as a "salad bowl" and as a "cultural mosaic", in contrast to a "melting pot". (Source: wikipedia, 2020)

Interculturalism refers to support for cross-cultural dialogue and challenging self-segregation tendencies within cultures. Interculturalism involves moving beyond mere passive acceptance of a multicultural fact of multiple cultures effectively existing in a society and instead promotes dialogue and interaction between cultures.

Interculturalism has arisen in response to criticisms of existing policies of multiculturalism, such as criticisms that such policies had failed to create inclusion of different cultures within society, but instead have divided society by legitimizing segregated separate communities that have isolated themselves and accentuated their specificity. It is based on the recognition of both differences and similarities between cultures. It has addressed the risk of the creation of absolute relativism within postmodernity and in multiculturalism. (Source: wikipedia, 2020)

Aguado proposes these principles (See “*La Educación Intercultural: Concepto, Paradigmas, Realizaciones*”, [Further reading], page 42.

1. Promote the respect by all cultures together and condemn the politics to change the culture of the people towards the culture dominant. (Borrelli y Essinger, 1989)
2. The intercultural education is relevant for any student, not only for the foreigners and minorities (Borrelli and Essinger, 1989)
3. The troubles created by the ethnic and cultural diversity of the society has many solutions, there not an only magic solution. The politics in education there are partials because we are in a global society (Galino, 1990).
4. It's based in the perception about to accept cultures in contact, it's near to the form of life of societies with a poor cultural context instead of societies with more rich, more structure and high social control.
5. We need develop a scheme of concepts with many cultures demonstrating in the education that the knowledge is the common property of all people (Walking, 1990).

So, interculturalism and multiculturalism are the same concept in many uses, both recognize the cultural diversity in the contexts where there are the diversity, but interculturalism is doing an emphasis in the enrichment of all cultures respecting the diversity.

Damegender understands has an international and intercultural perspective about guess the gender about the name in the sense that in many countries are existing many different cultures determining names, surnames with a gender. So, in Spain are living 4 so important cultures (no foreigners):

- Castilian (culture dominant)
- Catalan
- Basque
- Galician

These cultures has correlations with names and surnames.

7.3 Feminism, Ecofeminism and Intersectionality

Feminism is a range of social movements, political movements, and ideologies that aim to define and establish the political, economic, personal, and social equality of the sexes. Feminism incorporates the position that societies prioritize the male point of view, and that women are treated unjustly within those societies. Efforts to change that include fighting against gender stereotypes and establishing educational, professional, and interpersonal opportunities and outcomes for women that are equal to those for men. (Source: wikipedia, 2020).

Ecofeminism is a branch of feminism that sees environmentalism, and the relationship between women and the earth, as foundational to its analysis and practice. Ecofeminist thinkers draw on the concept of gender to analyse the relationships between humans and the natural world. The term was coined by the French writer Françoise d'Eaubonne in her book *Le Féminisme ou la Mort* (1974). Ecofeminist theory asserts a feminist perspective of Green politics that calls for an egalitarian, collaborative society in which there is no one dominant group. Today, there are several branches of ecofeminism, with varying approaches and analyses, including liberal ecofeminism, spiritual/cultural ecofeminism, and social/socialist ecofeminism (or materialist ecofeminism). Interpretations of ecofeminism and how it might be applied to social thought include ecofeminist art, social justice and political philosophy, religion, contemporary feminism, and poetry. (Source: wikipedia, 2020)

The goal 5 in United Nations in 2020 is “Achieve gender equality and empower all women and girls”. (Source: United Nations website)

Damegender don't reduce the gender gap per se. It's a tool to measure gender gap in Internet. The data is the basis to do politics to reduce the gender gap. These data must be used in contexts helping to the women, such as, feminism associations, political parties or trade unions with accomodation in favor to the gender equity, etc.

Damegender is giving more oportunities to the women to reduce the gender gap than another tools due to the license system.

Intersectionality is a theoretical framework for understanding how aspects of a person's social and political identities (e.g., gender, sex, race, class, sexuality, religion, disability, physical appearance, height, etc.) combine to create unique modes of discrimination and privilege. Intersectionality identifies advantages and disadvantages that are felt by people due to a combination of factors. For example, a black woman might face discrimination from a business that is not distinctly due to her race (because the business does not discriminate against black men) nor distinctly due to her gender (because the business does not discriminate against white women), but due to a unique combination of the two factors. (Source: wikipedia, 2020)

So, to understand discrimination, we must understand multiple factors. For example, the free software communities with the principles about no discrimination and share code seems a good place to advance in the rights of the women in the software world, but if the reality is a place dominated by men then we must look for another ideas. Intersectionality is about to find the best formula to advance in values understanding that the rights of the women depends of another values such as democracy, free speech, labor rights, ...

7.4 Gender

Gender is the range of characteristics pertaining to, and differentiating between, masculinity and femininity. Depending on the context, these characteristics may include biological sex, sex-based social structures (i.e., gender roles), or gender identity.

Most cultures use a gender binary, having two genders (boys/men and girls/women); those who exist outside these groups fall under the umbrella term non-binary or genderqueer. Some societies have specific genders besides "man" and "woman", such as the hijras of South Asia; these are often referred to as third genders (and fourth genders, etc.). (Source: wikipedia, 2020)

Transgender people have a gender identity or gender expression that differs from their sex assigned at birth. Some transgender people who desire medical assistance to transition from one sex to another identify as transsexual. Transgender, often shortened as trans, is also an umbrella term. In addition to including people whose gender identity is the opposite of their assigned sex (trans men and trans women), it may include people who are not exclusively masculine or feminine (people who are non-binary or genderqueer, including bigender, pangender, genderfluid, or agender). Other definitions of transgender also include people who belong to a third gender, or else conceptualize transgender people as a third gender. The term transgender may be defined very broadly to include cross-dressers. (Source: wikipedia, 2020)

In Damegender, we are applying binary ideology classifying people as male or female only due to that the free datasets provided by the states only supports this idea in the moment writing this book. But we respect the non binary philosophies due to this philosophies are describing a reality.

8 Conclusions

There are many options to count males and females in Internet, a good idea is to retrieve a dataset about males and females. Damegender is giving the most modern open datasets and it provides a good toolkit for many solutions:

- To count males and females in git repositories, mailing lists, csv files, ...
- To predict gender with machine learning if the name is not in the dataset
- To guess the country about the surname
- To understand how is used a name with different cultural regions
- To retrieve names from commercial apis
- To view relationships between names and races

These techniques can help to research or to visualize gender gap.

With this manual we have understood:

- How to use Damegender.
- How to compare different solutions with a scientific perspective, that is to manage mathematical vocabulary for to apply the concepts.
- How to apply this software to different use cases.
- How to find the main external resources about gender gap.
- To explain some philosophies for to give explanations to the data.

Acknowledgments

Damegender is dedicated to the INTER Project. I, David Arroyo Menéndez wants to give thanks to:

- To my mother.
- To Bitergia team by the motivation towards Python and Perceval
- To Jesús González Barahona and Gregorio Robles guiding me towards the success in the publication.
- To GAPLEN team (Leticia Martín Fuertes, Ana González Ledesma, ...) by the motivation towards Natural Language Processing.
- To Gema Rodríguez Pérez by the comments and motivation in the subject and the shared lectures.
- To the feminist movement by show me the problems related with the gender gap.
- To the Free Software movement by the understanding in licenses and the happy hacking sessions.
- To URJC by the seminars.
- To Damegender contributors of code (Luz Galvis, Oriol Tauferia and Jesús González Barahona)

Further reading

La Máquina Reaccionaria by María Ávila (Published by Tirant Humanidades)

The objective of this book is to analyze two social strengths. First, the fight of the women changing her position in the society and the structure of the patriarchy. Second, the resistances to decrease and to stop this change in a clear, systematic and aware form or in a silly and involuntary form.

Measuring the Gender Gap on the Internet by Bruce Bimber (Published by University of Texas Press)

This paper evaluates differences in men's and women's presence on the Internet, testing for the presence of gender-specific causes for different rates of Internet use. Methods. The paper presents new survey data collected by the author in 1996, 1998, and 1999 showing trends in Internet use, and presents regression models of Internet access and use. Results. Two statistically significant gender gaps exist on the Internet: in access and in use. The access gap is not the product of gender-specific factors, but is explained by socioeconomic and other differences between men and women. The use gap is the result of both socioeconomics and some combination of underlying gender-specific phenomena. Conclusions. Around one-half of the "digital divide" between men and women on the Internet is fundamentally gender related. Several possible causes may explain this phenomenon.

Women in STEM: A Gender Gap to Innovation by various authors (Published by Economics and Statistics Administration Issue Brief No. 04-11)

This executive summary very good referenced presents data about males and females in STEM. Very useful to compare data (for example, in this manual we have compared this data with males and females in the Free Software world).

Colapso. Capitalismo terminal. Transición Ecológica. Ecofascismo. by Carlos Taibo (Published by Catarata, ISBN 978-84-9097-203-8)

This book is a good explanation about politics chances related with ecology, climatic change, etc.

Free Software Free Society by Richard Stallman (Published by GNU Press, ISBN 1-882114-98-1)

Richard Stallman is the best philosopher about Free Software being the father of the most used free software licenses and the founder of GNU project being developed by him and many other people. This book explains all ethical ideas that is inspiring the free software community called GNU.

Comparison and benchmark of name-to-gender inference services by Lucía Santamaría and Helena Mihaljevic (Published by PeerJ Journal)

This paper has inspired the development of Damegender and many ideas about name-to-gender software is implemented in this book. Thanks a lot by this job.

The Effect of Gender in the Publication Patterns in Mathematics by Helena Mihaljevic, Lucía Santamaría, Marco Tullney (Published by Plos One Journal)

A very good scientific job about gender gap in science. This paper is suggested to people who has learnt to use Damegender and they want develop a good scientific job.

Damegender: Writing and Comparing Gender Detection Tools by David Arroyo Menéndez (Published by EasyChair 2020).

This paper presents the scientific perspective about Damegender.

Perceval: software project data at your will by Santiago Dueñas, Valerio Consentino, Gregorio Robles, Jesús M. González Barahona (Published by ICSE 2018).

This paper presents Perceval a software to retrieve information from different sources. Damegender is using perceval to retrieve data from git in the moment to write this book.

The gender gap in science: How long until women are equally represented? by Luke Holman, Devi Stuart-Fox, Cindy E. Hauser.

Using the PubMed and arXiv databases, the paper estimated the gender of 36 million authors from >100 countries publishing in >6000 journals, covering most STEMM disciplines over the last 15 years, and made a web app allowing easy access to the data (<https://lukeholman.github.io/genderGap/>).

La Educación Intercultural: Concepto, Paradigmas, Realizaciones by Teresa Aguado Odina

This paper is a very good reference to understand Interculturalism focused on education. The vocabulary explained in this paper can be so useful to apply in concepts related with feminism, or social sciences in general.

Guía INTER: una guía práctica para aplicar la educación intercultural en la escuela by Teresa Aguado Odina, Inés Gil Jaurena y otras personas. (Published by Universidad Nacional de Educación a Distancia UNED).

This document allows to apply the interculturalism concepts in the schools.

Machine Learning by Tom M. Mitchell (Published by Mc Graw Hill, ISBN: 0-07-042807-7)

Tom Mitchell is a father of the Machine Learning. This book explains in a simple way the main concepts understanding the need about the machine learning. So useful to learn algorithms.

Periodismo y Social Media: como estan usando Twitter los periodistas españoles by Clara Sainz de Baranda (Published by Estudios sobre el Mensaje Periodístico, UCM)

Understanding the Demographics of Twitter Users by Mislove, A., Lehmann, S., Ahn, Y. Y., Onnela, J. P., & Rosenquist, J. N. published by Icwsrm

This paper could have been written with Damegender. In this paper the author calculates gender of twitter users using first names and race/ethnicity using last names.

Discriminating Gender on Twitter by John D. Burger and John Henderson and George Kim and Guido Zarrella

This paper explains classify gender in Twitter obtaining very good results with tweet texts, screen name, description and full name. The best classifier was obtaining the 92% of accuracy.

Galaxia Internet by Manuel Castells (Published by Plaza & Janés, ISBN: 978-8401341571)

This book explains Internet with the point of view of a big sociologist. So it is regarded as a good introduction to Social informatics. That is the study of information and communication tools in cultural or institutional contexts.

Global Gender Gap Index 2012 (Published by World Economic Forum)

The Global Gender Gap Report was first published in 2006 by the World Economic Forum. The 2020 report (published in 2019) covers 153 countries. The Global Gender Gap Index is an index designed to measure gender equality.

Has Feminism Changed Science by Londa Schiebinger (Published by The University of Chicago Press Journals)

Feminism and Science by Evelyn Fox Keller & Helen E. Longino (Published by Oxford University)

Teoría Feminista by Ana de Miguel & Celia Amorós (Published by Minerva Ediciones)

Constructing Grounded Theory by Kathy Charmaz (Published by SAGE)

Grounded Theory is the most important theory to apply qualitative research. You can learn to classify and to count discourses (for example, feminist discourses) and to reach conclusions with this book.

Técnicas Cualitativas de Investigación Social by Miguel Valles (Published by Síntesis Sociología)

A good first book about qualitativism for sociology, politics or social work students, it's giving theoretical contents with practical examples done in Spain. If you want learn to put Damegender data into discourses this book can learn you to do it.

Encuesta sobre Equipamiento y Uso de Tecnologías de Información y Comunicación en los Hogares by INE.es (https://www.ine.es/prensa/tich_2020.pdf)

This document explains differences in gender using Internet and computers in Spain.

Reproducible Research in Computational Science by Roger D. Peng published by Science (<https://science.sciencemag.org/content/sci/334/6060/1226.full.pdf>)

This paper describes the quality of papers for reproducibility.

Appendix A License

Version 1.3, 3 November 2008

Copyright © 2000, 2001, 2002, 2007, 2008 Free Software Foundation, Inc.
<https://fsf.org/>

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

0. PREAMBLE

The purpose of this License is to make a manual, textbook, or other functional and useful document *free* in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or non-commercially. Secondarily, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of “copyleft”, which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

1. APPLICABILITY AND DEFINITIONS

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The “Document”, below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as “you”. You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A “Modified Version” of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A “Secondary Section” is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document’s overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The “Invariant Sections” are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released

under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The “Cover Texts” are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A “Transparent” copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not “Transparent” is called “Opaque”.

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTEX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The “Title Page” means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, “Title Page” means the text near the most prominent appearance of the work’s title, preceding the beginning of the body of the text.

The “publisher” means any person or entity that distributes copies of the Document to the public.

A section “Entitled XYZ” means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as “Acknowledgements”, “Dedications”, “Endorsements”, or “History”.) To “Preserve the Title” of such a section when you modify the Document means that it remains a section “Entitled XYZ” according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

2. VERBATIM COPYING

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

3. COPYING IN QUANTITY

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

4. MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any,

be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.

- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- D. Preserve all the copyright notices of the Document.
- E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- H. Include an unaltered copy of this License.
- I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
- J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
- L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.
- N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.
- O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their

titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

5. COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled "History" in the various original documents, forming one section Entitled "History"; likewise combine any sections Entitled "Acknowledgements", and any sections Entitled "Dedications". You must delete all sections Entitled "Endorsements."

6. COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

7. AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an “aggregate” if the copyright resulting from the compilation is not used to limit the legal rights of the compilation’s users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document’s Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

8. TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled “Acknowledgements”, “Dedications”, or “History”, the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

9. TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense, or distribute it is void, and will automatically terminate your rights under this License.

However, if you cease all violation of this License, then your license from a particular copyright holder is reinstated (a) provisionally, unless and until the copyright holder explicitly and finally terminates your license, and (b) permanently, if the copyright holder fails to notify you of the violation by some reasonable means prior to 60 days after the cessation.

Moreover, your license from a particular copyright holder is reinstated permanently if the copyright holder notifies you of the violation by some reasonable means, this is the first time you have received notice of violation of this License (for any work) from that copyright holder, and you cure the violation prior to 30 days after your receipt of the notice.

Termination of your rights under this section does not terminate the licenses of parties who have received copies or rights from you under this License. If your rights have been terminated and not permanently reinstated, receipt of a copy of some or all of the same material does not give you any rights to use it.

10. FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <https://www.gnu.org/licenses/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License “or any later version” applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation. If the Document specifies that a proxy can decide which future versions of this License can be used, that proxy’s public statement of acceptance of a version permanently authorizes you to choose that version for the Document.

11. RELICENSING

“Massive Multiauthor Collaboration Site” (or “MMC Site”) means any World Wide Web server that publishes copyrightable works and also provides prominent facilities for anybody to edit those works. A public wiki that anybody can edit is an example of such a server. A “Massive Multiauthor Collaboration” (or “MMC”) contained in the site means any set of copyrightable works thus published on the MMC site.

“CC-BY-SA” means the Creative Commons Attribution-Share Alike 3.0 license published by Creative Commons Corporation, a not-for-profit corporation with a principal place of business in San Francisco, California, as well as future copyleft versions of that license published by that same organization.

“Incorporate” means to publish or republish a Document, in whole or in part, as part of another Document.

An MMC is “eligible for relicensing” if it is licensed under this License, and if all works that were first published under this License somewhere other than this MMC, and subsequently incorporated in whole or in part into the MMC, (1) had no cover texts or invariant sections, and (2) were thus incorporated prior to November 1, 2008.

The operator of an MMC Site may republish an MMC contained in the site under CC-BY-SA on the same site at any time before August 1, 2009, provided the MMC is eligible for relicensing.

ADDENDUM: How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page:

Copyright (C) *year your name*.
Permission is granted to copy, distribute and/or modify this document
under the terms of the GNU Free Documentation License, Version 1.3
or any later version published by the Free Software Foundation;
with no Invariant Sections, no Front-Cover Texts, and no Back-Cover
Texts. A copy of the license is included in the section entitled ‘‘GNU
Free Documentation License’’.

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the “with... Texts.” line with this:

with the Invariant Sections being *list their titles*, with
the Front-Cover Texts being *list*, and with the Back-Cover Texts
being *list*.

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

Appendix B Photos

Photos taken in the Damegender development with a community giving mutual support.

The first photo is about Richard Stallman visiting Madrid in Fundación Telefonica. Thanks to Richard Stallman we have a philosophy and a free operating system where we can share code.

The second photo is about Jesús González Barahona, who were mentoring the work and giving good suggestion every week.

The third photo is about Gema Rodríguez Pérez celebrating the success Phd lecture with Gregorio Robles and another researchers and friends giving support to the research.

The fourth photo is about the Madrilenian Seminar on Software Development where the people was showing the advances in their works in research.









Index

A

Accuracy	10
Acknowledgments	41
APIs	1

B

Bibliography	42
--------------------	----

C

Choosing components	15
Commands	4
Commands about Statistics	4
Conclusions	40
Configuring API Keys	3
Confusion matrix	10
Counting features in names	14
Counting males and females in a git repository	25
Counting males and females in Debian	18
Counting males and females in Forbes	20
Counting males and females in Linux Kernel	20
Counting males and females in Maps	26

D

Damegender License	45
Debian	18
Deciding for males and females in images	22

E

Ecofeminism	33, 38
Error coded	10
Error coded without na	10
Error gender bias	10
Executing tests	4

F

F1 score	10
False negative	10
False positive	10
Feminism	33, 38
Forbes	18
Free Software	33
Further reading	42

G

Gender	28, 33, 39
Gender Detection Tools from the Name	1
Gender Gap	28
Gender gap in science	27
Gender Inequality in Free Software	31
Gender Inequality in STEM	30
Gender Inequality in the World	28
Git	18

I

Installation	3
Interculturalism	33, 36
Intersectionality	38
Introduction	1

K

Kernel	18
--------------	----

L

Las Calles de la Mujeres	18
--------------------------------	----

M

Maps	18
Measuring success and error	10
Multiculturalism	33, 36

O

Open Source	33
-------------------	----

P

Perceval	4
Precision	10
Principal Component Analysis (PCA)	10, 14
Python Virtual Environment	3

Q

Qualitative Research	33
----------------------------	----

R

Recall	10
Regenerating files in post installation	4
Reproducible Research	1
Richard Stallman	33
ROC	10

Index	58
-------	----

S

- Secondary Sources 28
- Social Research 28, 33
- Software Philosophies 33

T

- Theoretical Frameworks 33
- True negative 10
- True positive 10

W

- Webscraping 18
- Webscraping and Damegender
(counting scholars) 22