

Damegender: Towards an International and Free Dataset about Name, Gender and Frequency

David Arroyo Menéndez
Grupo de Sistemas y Comunicaciones (GSyC)
Universidad Rey Juan Carlos, Madrid, Spain
{d.arroyome@alumnos}.urjc.es

URJC

Abstract

Equality of gender is the 5th objective of sustainable development in United Nations¹.

This equality can be reached working on to measure and to analyze data and to apply politics from the results. On many gender studies, we need to count males and females deciding gender from names, for instance, research papers, job positions, streets, ... The traditional way is to use commercial APIs with proprietary data without idea about how the data has been built. Another way, is taking data from Wikipedia, linguistic studies, scientific sites, ...

Many statistical institutions are providing Open Datasets about name, gender and frequency. So, we need a scientific discussion about unifying formats, making easy ways to process these data and ways towards make standards.

Meanwhile, has been developed Damegender (Free and Open Source Software) to retrieve and make calculus with these data.

The dataset is covering more than 20 countries in the occidental world reaching a big number of names and accuracies around (87.56%) with it. Surely, allowing to measure gender gap to students and academics interested on the phenomenon without costs and on a reproducible

way, more people will be contributing to fix the gender gap.

There are a warranty of quality on reproducible research, that's the Free Software and the citation about official sources about names, gender and frequency provided by statistics institutions making easy the peer review and opening doors to the semantic web and the attention to diversity.

1 Introduction

There are a goal in United Nations about gender gap². What can I do as software engineer?. The first step is to remember that “if you cannot measure it, you cannot improve it” [Tho33] and “Software Engineering Economics is an invaluable guide to determining software costs, applying the fundamental concepts of microeconomics to software engineering [B⁺81]”. A common profit using Free Software and Open Data used to be the reduction of costs, for example, many people and institutions is using LibreOffice and Ubuntu (GNU/Linux) to avoid the payment by proprietary licenses with similar solutions such as Microsoft Windows and Microsoft Office, the new costs used to be change a social inertia changing popular products by products without costs. There are a set of gender detection tools from the name is based on API solutions, giving a free software and open data solution, it will be making a competition in a market without a very strong leader, avoiding payments and thinking on strategies about profits from a trademark, such as, Firefox or Chrome.

Detecting personal names, you can infer gender on academical papers, books, newspapers and many interactions on Internet. So, the task about detect gender

Copyright © by the paper's authors. Use permitted under Creative Commons License Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)

¹<https://www.un.org/sustainabledevelopment/gender-equality/>

²<https://www.un.org/sustainabledevelopment/gender-equality/>

from the names could be a strategic task to measure gender gap.

Nowadays, many people are using APIs such as Genderapi, Genderize, Namsor, or NameApi. Another people is using solutions based on Wikipedia, or Free Software solutions (NLTK[LB02], R Gender, Gender Detector, Gender Computer³, ...)

Open Source solutions has a few number of names due to use files of a single country or being software not maintained in the long time. And Wikipedia is not taking into account the frequency of the personal names.

However, the gender gap is a problem recognised in United Nations and the IT market is leading big inequalities in the world in economy and gender gap. This paper presents a real work collecting data with a scientific perspective to solve the problem giving a full toolkit that is solving a good number of problems (search engine, infer gender in csv files, names in different countries, wide dataset, ...) faced by the industry and other problems not solved in an industrial way such as count males and females in Github repositories, mailing lists, ...

Another previous work [KWL⁺16] about this kind of tools is discussing about the datasets as a way to improve the accuracies, comparing tools that is using different public datasets (SSA, IPUMS, namdict, ...)

The solutions is being faced by a practical way: augmenting the number of names using official statistics and taking into account diversity goals such as non binary gender and cultural minorities.

With DameGender we will make science reproducible[Pen11] in fields where has done similar works such as Natural Language Processing (gender detection from the name [SGT⁺19]), social sciences or journalism (gender gap [HSFH18, MLA⁺11, NP17, dBA14]), linguistic [Hut16, vdWRvdW⁺20, Oka18], software engineering [VCS12], ...

The remainder of this paper is structured as follows:

In Section 2 presents the main research works about to measure gender gap and gender detection tools from the name.

Section 3 is giving vocabulary and philosophy about to choose sources and to face the diversity troubles building a dataset.

Section 4 explains an application about this dataset to measure gender gap in GNU/Linux.

Section 5 points a summary about this approach and future works.

The contributions to the State of Art presented in this paper are:

1. An integrated solution where make experiments

³<https://github.com/tue-mdse/genderComputer>

in the different applications field relative to inferring gender from the name.

2. A collection of Open Datasets retrieved from statistical sources and standardized in an unique format about gender, name and frequency

3. A new study applying DameGender to count males and females in GNU/Linux

4. An approach based on reproducible results.

2 State of Art

2.1 About Gender Gap

Fix the gender gap refers to equality between males and females, it's about non discrimination policies between males and females. The gender is about the sex determined in the moment of the birth, although can be changed in some moment in the life. So, there are discussions about the gender definition referring to these problems. But there are consensus determining gender, frequency and names with official statistics released by the institutions in the states.

To measure gender gap requires set indicators about it. In [For21] has been proposed economy, health, education and politics. And in United Nations⁴ there are indicators such as laws, education, mortality of mothers, political participation, poverty, domestic work, gender parity in the work, access to economy, situation about the youth (access to studies and/or work), violence against the women, climate justice, access to the justice, health, ...

When the measurement has been done with these indicators defined, it's possible make decisions improving the situation through conclusions about researches, for example, in [MKSF⁺10] are concluding that making affirmations about ethical values are being reduced the gender achievement gap in colleges.

Many times, to measure the gender gap has been reached with methodologies in social research, such as the survey. For example, in [Bim00] has been presented two factors in the gender gap in Internet (access and use) by socioeconomic and gender reasons in a survey collecting data for several years. The Internet access is so important in the economy, education, ... Another work [RRGBD16] is using a survey of 2000 contributors.

2.2 Counting Males and Females in Internet. Why? Where?

This work is focused retrieving data from secondary sources such as GitHub, Wikipedia, APIs, websites in general, mailing lists, ... There are serious previous research works about factors modifying several gender

⁴<https://www.unwomen.org/>

gap indicators (economy, education, politics, ...) from secondary sources.

For example, a social scientist studying gender gap in journalism [ÁACS12] can be counting males and females in Twitter. These metrics are important because the journalism is determining gender gap or not in political, education, economy, ... Meanwhile, there are Computer Science making research about how to count males and females in Twitter [BHKZ11]. On these works, is being retrieved name, nickname, photo and identifying gender from these data.

[BHKZ11] has been presenting several configurations of a language-independent classifier for predicting the gender of Twitter users. The large dataset used for construction and evaluation of these classifiers was drawn from Twitter users who also completed blog profile pages.

And in [MLA⁺11] has been analyzed the Twitter population including the gender. The gender was inferred making queries from the names to the dataset provided by US statistic institution.

In [WGJS15] has been analyzing the gender gap in Wikipedia showing evidence of more subtle forms of gender inequality, being Wikipedia an interesting factor in education, explaining how to solve these evidences. To measure gender inequality has been developed the next bias: coverage, structural, lexical (ex: discriminatory words for women) and visibility.

Computer Science is generating many Forbes billionaires, the public code is a factor to understand the gender gap about Computer Science and that's has some importance in economy. The public repositories can be used to build indicators about economy in Computer Science with more factors, such as job positions, value of companies, ... In [Zac20] has been conducted the first large-scale longitudinal study of gender imbalance among authors of collaboratively developed, publicly available code, where contributions by female authors remain scarce: less than 8 % of commits for which we could detect a gender, confirming decades of gender imbalance in Free/Open Source Software (FOSS). Steffano was using namdict dataset through of gender-guesser to infer gender from the name. Another work related is [VPR⁺15] explaining that women programmers are in the minority in OSS and other technical teams, although increased gender and tenure diversity are associated with greater productivity.

In [VCS12] are exploring the popular Q&A about technological stuff called StackOverflow summarizes that the percentage of women engaged in SO is greatly imbalanced, and men represent the vast majority of contributors.

And [IHSR18] are showing data few females participating contributing code and taking political responsibilities in OpenStack community.

Related to gender gap in science is so interesting [HSFH18] presenting a code in R where making calls to genderize api is giving good approaches about how to calculate gender gap inferring gender from authors name retrieved from arXiv.

2.3 Automatic approaches to infer gender

There are several tasks related with infer gender from Internet sources: hand written, images, documents, names, ...

In [LSB11] presents a method inferring gender from hand written texts with a 67.5 % of accuracy.

Related to infer gender from images [GC08] combines image based gender and age classifiers with the cultural context information provided by first names to recognize people with no labeled examples with results near to 60 % of accuracy.

In [AKFS03] explains that the females uses many more pronouns and males use many more noun specifiers in a large subset of the British National Corpus covering a range of genres. So, in [KAS02] is being presented a document classification system with accuracy of approximately 80 per cent. And in [CCS11] exposes a feature selection and a model built using Machine Learning earning an accuracy of 85.1 per cent identifying gender from text.

2.4 Inferring gender from the name

Generally, the tools to infer gender from the name is based on datasets that is including gender and name as minimum.

In [LR13] has been presented a method to infer gender from first names in Twitter, the dataset was built on hand coded by the agreement of 3 Amazon workers with 50000 Twitter users select at random out with only 12681 gender labels. The goal of this study was to determine the incremental value of using the user name as a feature in gender inference based on tweets.

In [MS16] presents how to infer gender in Twitter. They are using namdict and us census as datasets. The features are 'number of consonants', 'number of vowels', 'number of syllables', 'number of bouba consonants', 'number of bouba vowels', 'number of kiki consonants', 'number of kiki vowels'. The classification model is using SVM.

2.5 Other ideas related

So, in [AWM⁺09] was presented a system to classify name and ethnicity from Open Sources using Machine Learning extracting a name list from Wikipedia. A more recent work is [NRN21] where is presented Nam-Prism being applied to massive software repositories.

Another approach was presented in [BMI10] a lexical-pattern-based approach to extract aliases of a given name. With a set of names and their aliases as training data to extract lexical patterns. The candidates are ranked using various ranking scores. And to construct the ranking function was used ranking Support Vector Machines.

2.6 Standards related

ISO/IEC 5218 proposes a norm about coding gender: “0 as not know”, “1 as male”, “2 as female” and “9 as not applicable”.

The RFC 6350 (vCard) ⁵ where the section Gender has these categories: “m as male”, “f as female”, “o as other”, “n as not applicable” and “u as undefined”. Based on this standard the people who is doing web publishing can use css classes using a web standard such a h-card ⁶ microformats. And in the context of to write forms in web interfaces consider w3 lectures ⁷

2.7 Summary

In the State of Art of gender inference is the first name the key factor to determine gender, although in many contexts there are more features: surnames, text, images, nicknames, ... The first name can be useful to infer another stuff such as race, ethnicity or culture, too.

Machine Learning and the previous features selection is being used in many works, although the discussion about what is the best approach is an open discussion.

The datasets can be built on hand by human experts, although there are some Open Datasets used several times in these researches, such as namdict, or us census.

3 Design

3.1 Truth and Falsehood in names, gender and frequency

The current idea in the field is the data about name, gender and frequency is ok because there are people who is paying by it, or many people is downloading a product. This intuition is right generally, although sometimes the people is paying by a bad product due to a good marketing strategy, a monopoly or there is a fraud, ... Another idea is the people trust in the government about statistics such as economy, demography, democracy, ... So the people can trust on names,

gender and frequency. In Damegender, we are trusting in both notions about truth: the market’s point of view and the official statistic’s point of view.

Sometimes there are problems downloading the official statistics, but there are people who has retrieved these data, for example, with web scraping. We want classify these files with another idea about truth.

Another problem arises when the government does little chances in the data, sometimes communicating it to the users and other times not. That could be a problem about upgrades, but it’s not a problem with the truth, although it’s possible make a trace about these chances.

Other sources to retrieve gender and names can be personal scientific websites, Wikipedia, or similar, but these sources is not giving the frequency, now. So, we are rejecting this idea.

With an international free dataset about names, gender and frequency we can build reproducible science in fields such as Natural Language Processing (gender detection from the name), social sciences or journalism (gender gap [HSFH18, MLA⁺11, NP17, dBA14]), linguistic [LN05, Kru62, vdWRvdW⁺20, Agy06, FMO⁺87], software engineering [VCS12], ...

3.2 Gender, Language, Nation and Diversity

There are rules and exceptions in the languages to predict if a name is about male or female when you don’t know the name. For example, in Spanish or English there are more names ending with ‘a’ classified as females than classified as males. And Andrea is female in Spain and male in Italy. So, it’s useful to understand the language and culture associated with a name. Language is close to nation, but there are differences, for example, in Spain there are several languages basque, catalan, castillian, ... or the Spanish is the main language in Spain and in other countries such as Argentina, Mexico, Ecuador, Bolivia, ... So, it would be useful to detect the language and nation from names and surnames to help to detect gender.

Some countries, such as Spain, are providing free datasets about surnames but we need more efforts from many countries on this objective. On other hand, there are previous works to relate name and surnames with ethnicity using Wikipedia and Machine Learning [AWM⁺09].

3.3 Damegender Open Datasets Collection

In Damegender, we have unified the different formats to name, gender and frequency from official sources in these countries: Argentina, Austria, Australia, Belgium, Canada, Denmark, Germany, Spain, Finland, France, Great Britain, Ireland, Mexico, New Zealand,

⁵<https://datatracker.ietf.org/doc/html/rfc6350>

⁶<https://github.com/microformats/h-card>

⁷<https://www.w3.org/International/questions/qa-personal-names>

Dataset	SSA	namdict	NLTK	Damegender
males	91.320	48.821	2.943	278.928
females	91.320	48.821	5.001	299.870

Table 1: Comparison about the number of names between Open Data solutions

Norway, Russia, Portugal, Slovenia and United States of America.

We have found 2 main criteria counting males and females: number of births in a year and people using the name in a year. So, we have divided the files being to able of make merging. The criteria making the count is: the average of the last twenty years where the data has been provided. Both criteria are good indicators to understand how many people is using a name as male or as female.

Later, we have merged these datasets building a free and international dataset.

We have found open datasets about countries such as Turkey and China retrieved by other open source developers that is being included in Damegender, but not in the international dataset. In Turkey the data has been retrieved using web scraping. And in China the data has been built by a company in collaboration with the China government and contributed to R language program. We want compare precision about this dataset with the commercial solutions to understand the truth about these datasets.

We have found surnames given by statistical institutions in Spain, Russia, United States of America and Argentina. So few statistical institutions is giving surnames. Understanding the diversity problem with this fact, we have retrieved surnames for all countries from Wikidata, these datasets contains few surnames being compared with datasets provided by statistical institutions. Although in names the diversity problem would be a minor problem we are giving names retrieved with Wikidata to the public.

When the work is finished, we could to rebuild machine learning models to predict new names and nicknames in any language and culture. The results is the longest list of public names.

A possible criticism about our idea is the Leslie Problem[BM15]: the match between gender and name has been changing in some years. And the answer is about you need introduce the age of the person to solve it. The most used use case is the input is the name and the output must be gender, frequency and percentage. So, we are deciding without age, surname, ... in the most of use cases. The idea about this dataset is to be designed for the most used use case. Although, we can take into account other inputs, such as surname or age to improve the accuracy. There are many Open Datasets with names and frequencies classified by years. So, this problem can be fixed with Open

Dataset	Accuracy	Precision	Recall	F1-Score
Damegender	0.8756	0.9638	1.0	0.925

Table 2: Several precision measures about the Damegender International Dataset

Data, too.

We have made measurements about the international DameGender dataset, using as base of truth the dataset explained in [SM18] reaching accuracy (0.8756), precision (0.9638), recall (1.0) and f1-score (0.925).

3.4 Free APIs for Free Datasets?

Many websites about Open Data is delivering methods to retrieve Open Data with a structured format and without costs for the user, such as, Wikipedia with SPARQL, OpenStreetMap with API rest, ...

We are detecting that the Open Datasets about names, gender and frequency are being modified one time per year as maximum for each statistical institution.

Damegender contains python scripts designed to create the different datasets and publish json files that could be used as a Free API Rest publishing the json files in sites as github pages, gitlab pages, or similar sites with free uploads.

```
$ cat DAVID_all.json
[{
  "name": "DAVID",
  "frequency": 4856689,
  "males": "99.73267796229078 %",
  "females": "0.26732203770922947 %"
}]
```

So, we could think that could have Free API Rest about names, gender and frequency reaching reduce costs to fix the gender gap on a collaborative way similar to Wikipedia, OpenStreetMap or many Free Software projects.

4 Measuring Gender Gap. GNU/Linux as Use Case

With a trust open dataset about names, gender and frequency is too easy to measure gender gap. Doing cheap to measure gender gap more students and academic people could work in the fifth Objective Development Sustainable of United Nations: to delete the gender gap.

This section is divided counting males and females in Debian, GNU and Linux.

We have reached the csv files from different ways to know the names about the people in these communities.

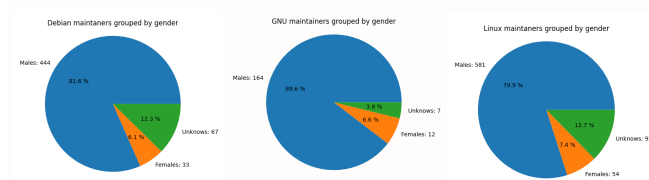


Figure 1: Males (blue), Females (orange) and Unknowns (green) in Debian, GNU and Linux

When this paper was being wrote in the Debian community all members must be collaborating with a gpg key, so we can count males and females from the keyring. The keyring was imported with gpg commands and later was dumped the keyring in a csv file.

In the moment to write this paper GNU⁸ and Linux⁹ has websites with the people collaborating in these projects. So, making web scraping scripts we have downloaded the people and processed the people to csv files

In Damegender, we have developed csv2gender, a software with a csv file as input and deploy a statistics graph and/or return the result of males, females and unknowns about the input.

To make easy to reproduce the experiment we are pasting the commands used with the version 0.3.4 of Damegender.

```
python3 csv2gender.py files/gnu.csv
--first_name_position=0
--title="GNU maintainers grouped by gender"
--dataset="inter"
--outcsv="files/gnu.gender.csv"
--outimg="files/gnu.gender.png"
--noshow --delete_duplicated

python3 csv2gender.py files/linux.csv
--first_name_position=0
--title="Linux maintainers grouped by gender"
--dataset="inter"
--outcsv="files/linux.gender.csv"
--outimg="files/linux.gender.png"
--noshow --delete_duplicated

python3 csv2gender.py files/debian.csv
--first_name_position=0
--title="Debian maintainers grouped by gender"
--dataset="inter"
--outcsv="files/debian.gender.csv"
--outimg="files/debian.gender.png"
--noshow --delete_duplicated
```

The inter dataset was created merging several open datasets downloaded from official statistics sites

⁸<https://www.gnu.org/people/>

⁹<https://www.kernel.org/doc/html/latest/process/maintainers>

from different nations: Austria, Australia, Belgium, Canada, Germany, Denmark, Spain, Finland, Ireland, Iceland, Mexico, New Zealand, Portugal, Slovenia, United States of America, Uruguay and France. That's a good representation of the Western World and the Free Software world is populating this world's area[GBRAIG08].

Linux divides the developers in 537 males (73.9%), 98 females (13.5%) and 92 unknowns (12.7%). The number of unknowns is due to different reasons, but it's so common in Linux that the developer is a company and not a name of a person.

GNU divides the developers in 164 males (89.6%), 12 females (6.6%) and 7 unknowns (3.8%)

The GNU people has a number lowest in females, they are the founder of the Free Software philosophy, the Debian principles and the Open Source philosophy was invented later influenced by GNU with very similar practical decisions (for example: deciding licenses for the software). Richard Stallman returned to be president recently apologizing by his personal behaviour with the females.¹⁰

Debian is a distribution, the project who makes the CD/DVD and the software ready to be downloaded from Internet with the dependencies. There are many distributions, such as, Ubuntu or RedHat so it is not representative, but it's interesting to understand that the numbers are similar in Debian dividing the developers in 408 males (75%), 69 females (12.7%) and 67 unknowns (12.3%).

5 Conclusions and Future Works

This paper is explaining the application about Damegender, the motivations (reproducible research, fix gender gap to reach an objective of United Nations, fields of application: linguistic, social sciences, software engineering, natural language processing, journalism, ...)

A good improvement is to build an international, universal and free dataset about names, gender and frequency with the right design with the current state of the job, attending to the diversity (LGBT options, cultural minorities, ...).

This paper has explained what technologies is involved on reduce costs about gender gap (gender detection from the names, api rest, semantic web, ...)

Augmenting the number of countries with statistical institutions giving names, gender and frequencies with Open Data will be augmenting the accuracies and giving more attention to the diversity.

The current state of work is the longest Open Dataset about names, gender and frequency with more

¹⁰<https://www.fsf.org/news/rms-addresses-the-free-software-community>

than 20 countries representing the Western World, being a solution with low number of unknowns in the real world.

The future works is about changes in the big software industry.

Making searches with strings about personal names (ex: Leticia) in search engines such as Google, these strings are not being classified as personal names, one solution will be data structured such as JSON-LD, microdata, microformats, rdfa ... Another solution will be store in the servers the Open Data Collection about names, gender and frequency and identify the context about the string is a personal name, that's an easy problem in popular sites such as Wikipedia, academic websites, ...

If the search engine identify the string as personal name, it can help to the user about the gender. That is similar than other problems such as streets, products, ... where you are giving additional information such as maps in streets or prices in products.

Other sites such as Github or Gitlab could be giving data about gender of developers in the site or in the software project with these datasets.

Another industry is about match sites (Meetic, Tinder, ...) where only is important photos, age and gender generally. It could be possible to give to users gender, photo and interests from personal names our open data collection and information related in Internet.

References

- [ÁACS12] Pilar Carrera Álvarez, Clara Sainz De Baranda Andújar, Eva Herrero Curiel, and Nieves Limón Serrano. Journalism and social media: How spanish journalists are using twitter/periodismo y social media: cómo están usando twitter los periodistas españoles. *Estudios sobre el mensaje periodístico*, 18(1):31, 2012.
- [Agy06] Kofi Agyekum. The sociolinguistic of akan personal names. *Nordic journal of African studies*, 15(2):206–235, 2006.
- [AKFS03] Shlomo Argamon, Moshe Koppel, Jonathan Fine, and Anat Rachel Shimoni. Gender, genre, and writing style in formal written texts. *Text & Talk*, 23(3):321–346, 2003.
- [AWM+09] Anurag Ambekar, Charles Ward, Jahangir Mohammed, Swapna Male, and Steven Skiena. Name-ethnicity classification from open sources. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge Discovery and Data Mining*, pages 49–58, 2009.
- [B⁺81] Boehm Barry et al. Software engineering economics. *New York*, 197, 1981.
- [BHKZ11] John D Burger, John Henderson, George Kim, and Guido Zarrella. Discriminating gender on twitter. In *Proceedings of the conference on empirical methods in natural language processing*, pages 1301–1309. Association for Computational Linguistics, 2011.
- [Bim00] Bruce Bimber. Measuring the gender gap on the internet. *Social science quarterly*, pages 868–876, 2000.
- [BM15] Cameron Blevins and Lincoln Mullen. Jane, john... leslie? a historical method for algorithmic gender prediction. *DHQ: Digital Humanities Quarterly*, 9(3), 2015.
- [BMI10] Danushka Bollegala, Yutaka Matsuo, and Mitsuru Ishizuka. Automatic discovery of personal name aliases from the web. *IEEE Transactions on Knowledge and Data Engineering*, 23(6):831–844, 2010.
- [CCS11] Na Cheng, Rajarathnam Chandramouli, and KP Subbalakshmi. Author gender identification from text. *Digital Investigation*, 8(1):78–88, 2011.
- [dBA14] Clara Sainz de Baranda Andújar. El género de los protagonistas en la información deportiva (1979-2010): noticias y titulares/the gender of the main characters in sports reporting (1979-2010): News and headlines. *Estudios sobre el mensaje periodístico*, 20(2):1225, 2014.
- [FMO+87] Peter Marshall Fraser, Elaine Matthews, Michael J Osborne, Sean G Byrne, Richard WV Catling, J-S Balzat, E Chiricat, Thomas Corsten, and Fabienne Marchand. *A lexicon of Greek personal names*, volume 5. Lexicon of Greek Personal Name, 1987.

- [For21] World Economic Forum. Global gender gap report 2021: Insight report. World Economic Forum Cologny, Switzerland, 2021.
- [GBRAIG08] Jesus M Gonzalez-Barahona, Gregorio Robles, Roberto Andradas-Izquierdo, and Rishab Aiyer Ghosh. Geographic origin of libre software developers. *Information Economics and Policy*, 20(4):356–363, 2008.
- [GC08] Andrew C Gallagher and Tsuhan Chen. Estimating age, gender, and identity using first name priors. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008.
- [HSFH18] Luke Holman, Devi Stuart-Fox, and Cindy E Hauser. The gender gap in science: How long until women are equally represented? *PLoS biology*, 16(4):e2004956, 2018.
- [Hut16] Matthew Hutson. The gender of names. *Scientific American Mind*, 27(4):14–14, 2016.
- [IHSR18] Daniel Izquierdo, Nicole Huesman, Alexander Serebrenik, and Gregorio Robles. Openstack gender diversity report. *IEEE Software*, 36(1):28–33, 2018.
- [KAS02] Moshe Koppel, Shlomo Argamon, and Anat Rachel Shimoni. Automatically categorizing written texts by author gender. *Literary and linguistic computing*, 17(4):401–412, 2002.
- [Kru62] John R Krueger. Mongolian personal names. *Names*, 10(2):81–86, 1962.
- [KWL⁺16] Fariba Karimi, Claudia Wagner, Florian Lemmerich, Mohsen Jadidi, and Markus Strohmaier. Inferring gender from names on the web: A comparative evaluation of gender detection methods. In *Proceedings of the 25th International conference companion on World Wide Web*, pages 53–54, 2016.
- [LB02] Edward Loper and Steven Bird. Nltk: the natural language toolkit. *arXiv preprint cs/0205028*, 2002.
- [LN05] Edwin D Lawson and Natan Nevo. Russian given names: Their pronunciation, meaning, and frequency. *Names*, 53(1-2):49–77, 2005.
- [LR13] Wendy Liu and Derek Ruths. What’s in a name? using first names as features for gender inference in twitter. In *2013 AAAI Spring Symposium Series*, 2013.
- [LSB11] Marcus Liwicki, Andreas Schlapbach, and Horst Bunke. Automatic gender detection using on-line and off-line information. *Pattern Analysis and Applications*, 14(1):87–92, 2011.
- [MKSF⁺10] Akira Miyake, Lauren E Kost-Smith, Noah D Finkelstein, Steven J Pollock, Geoffrey L Cohen, and Tiffany A Ito. Reducing the gender achievement gap in college science: A classroom study of values affirmation. *Science*, 330(6008):1234–1237, 2010.
- [MLA⁺11] Alan Mislove, Sune Lehmann, Yong-Yeol Ahn, Jukka-Pekka Onnela, and J Niels Rosenquist. Understanding the demographics of twitter users. In *Fifth international AAAI conference on weblogs and social media*, 2011.
- [MS16] Juergen Mueller and Gerd Stumme. Gender inference using statistical name characteristics in twitter. In *Proceedings of the The 3rd Multidisciplinary International Social Networks Conference on Social Informatics 2016, Data Science 2016*, pages 1–8, 2016.
- [NP17] Mari K Niemi and Ville Pitkänen. Gendered use of experts in the media: Analysis of the gender gap in finnish news journalism. *Public Understanding of Science*, 26(3):355–368, 2017.
- [NRN21] Reza Nadri, Gema Rodriguezperez, and Meiyappan Nagappan. On the relationship between the developer’s perceptible race and ethnicity and the evaluation of contributions in oss. *IEEE Transactions on Software Engineering*, 2021.

- [Oka18] Benard Odoyo Okal. A linguistic overview of the patronymic and gender names amongst the selected african communities. 2018.
- [Pen11] Roger D Peng. Reproducible research in computational science. *Science*, 334(6060):1226–1227, 2011.
- [RRGBD16] Gregorio Robles, Laura Arjona Reina, Jesús M. González-Barahona, and Santiago Dueñas Domínguez. Women in free/libre/open source software: The situation in the 2010s. In Kevin Crowston, Imed Hammouda, Björn Lundell, Gregorio Robles, Jonas Gamalielsson, and Juho Lindman, editors, *Open Source Systems: Integrating Communities*, pages 163–173, Cham, 2016. Springer International Publishing.
- [SGT⁺19] Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElShrief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. Mitigating gender bias in natural language processing: Literature review. *arXiv preprint arXiv:1906.08976*, 2019.
- [SM18] Lucía Santamaría and Helena Mihaljević. Comparison and benchmark of name-to-gender inference services. *PeerJ Computer Science*, 4:e156, July 2018.
- [Tho33] W Thompson. Electrical units of measurement, popular lectures, 1833.
- [VCS12] Bogdan Vasilescu, Andrea Capiluppi, and Alexander Serebrenik. Gender, representation and online participation: A quantitative study of stackoverflow. In *2012 International Conference on Social Informatics*, pages 332–338. IEEE, 2012.
- [vdWRvdW⁺20] Jeroen van de Weijer, Guangyuan Ren, Joost van de Weijer, Weiyun Wei, and Yumeng Wang. Gender identification in chinese names. *Lingua*, 234:102759, 2020.
- [VPR⁺15] Bogdan Vasilescu, Daryl Posnett, Baishakhi Ray, Mark GJ van den Brand, Alexander Serebrenik, Premkumar Devanbu, and Vladimir Filkov. Gender and tenure diversity in github teams. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*, pages 3789–3798. ACM, 2015.
- [WGJS15] Claudia Wagner, David Garcia, Mohsen Jadidi, and Markus Strohmaier. It’s a man’s wikipedia? assessing gender inequality in an online encyclopedia. In *Ninth international AAAI conference on web and social media*, 2015.
- [Zac20] Stefano Zacchiroli. Gender differences in public code contributions: a 50-year perspective. *IEEE Software*, 38(2):45–50, 2020.