



“Desambiguando” el género: cómo obtener género de nombres de manera sistemática, gratuita y consciente

Elvira González-Salmón

#YoSigoPublicando

28 noviembre 2023

¿A qué nos referimos con desambiguación* de género?

Lista de nombres



"Algoritmo"



Lista de nombres y género



Estudios que tienen en cuenta el género

Ejemplo: bibliometría

PSYCHOLOGY & SEXUALITY
2021, VOL. 12, NO. 4, 332-344
<https://doi.org/10.1080/19419899.2020.1729844>

Routledge
Taylor & Francis Group

OPEN ACCESS

What is gender, anyway: a review of the options for operationalising gender

Anna Lindqvist ^{a,b}, Marie Gustafsson Sendén^b and Emma A. Renström^c

^aDepartment of Psychology, Lund University, Lund, Sweden; ^bDepartment of Psychology, Stockholm University, Stockholm, Sweden; ^cDepartment of Psychology, University of Gothenburg, Gothenburg, Sweden

ABSTRACT
In the social sciences, many quantitative research findings as well as presentations of demographics are related to participants' gender. Most often, gender is represented by a dichotomous variable with the possible responses of woman/man or female/male, although gender is not a binary variable. It is, however, rarely defined what is meant by gender. In this article, we deconstruct the concept 'gender' as consisting of several facets, and argue that the researcher needs to identify relevant aspects of gender in relation to their research question. We make a thorough exposition of considerations that the researcher should bear in mind when formulating questions about each facet, in order to exemplify how complex this construct is. We also remind the researcher that gender is not a binary category and discuss challenges in the balance between taking existing gender diversity into account and yet sorting participants into gender categorisations that function in statistical analyses. To aid in this process, we provide an empirical example on how gender identity may be categorised when using a free-text response. Lastly, we suggest that other measurements than participants' gender might be better predictors of the outcome variable.

ARTICLE HISTORY
Received 1 March 2019
Accepted 10 February 2020

KEYWORDS
Gender; gender identity; transgender; research methods; cisnormativity

* O estimación, identificación, asignación, etc.

Organización del curso

1. Utilidad del curso
2. Cuestiones éticas
3. Funcionamiento básico
4. Cuestiones a considerar
5. Opciones
6. Casos prácticos
7. Bibliografía



Ej. NamSor

1. Utilidad del curso

1. Utilidad del curso: ¿Cuándo identificar el género de nombres de manera sistemática?

1. Número considerable de datos
2. No hay opción de autoidentificación
 - Encuesta
3. Es el único dato que tenemos (i.e. [Face++](#), pronombres)
4. Cuando tiene sentido para nuestra investigación ([Lindqvist et al., 2021](#))
 - ¿Queremos saber el género o aspectos físicos o por ejemplo género legal?
 - ¿Es realmente parte de la investigación?



Conclusión: evitarlo siempre que se pueda



2. Cuestiones éticas

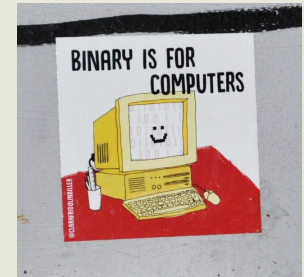
2. Cuestiones éticas

- Vemos la **probabilidad** de que un nombre sea asignado a un hombre o mujer, no asignamos género a personas individuales
- Habrá errores en esa probabilidad (Misgendering)
- Asume género de nombres - Refuerzo de estereotipos



2. Cuestiones éticas

- Imagen binaria del género (Frohard-Dourlent et al., 2016; Medeiros et al., 2020)
 - [Understanding Nonbinary People: How to Be Respectful](#)
 - [¿Qué es el género no binario?](#)
- No sabemos cómo de grande es el error

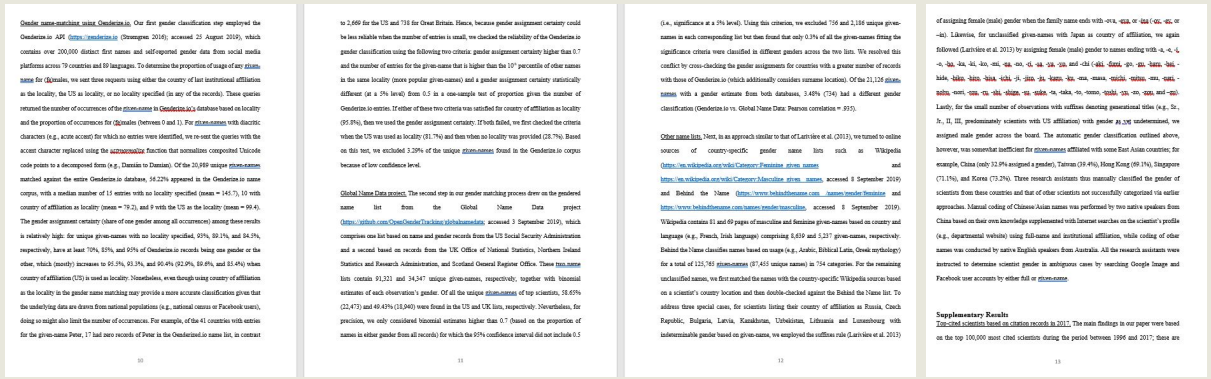


While the data collection for She Figures only considers sex-disaggregated data for men and women, it will be important to also consider non-binary gender for data collection in future publications, where possible. Non-binary is an umbrella term for gender identities that fall outside the gender binary of man or woman. This includes individuals whose gender identity is neither exclusively man nor woman, a combination of man and woman or between or beyond genders. The United Nations Economic Commission for Europe (UNECE)

[She Figures, 2021](#)

2. Cuestiones éticas

- Mucha gente trabajando en esto
- Importancia de los *disclaimers* y reconocimiento de limitaciones
 - Cada vez hay más
 - De “simply by looking at the full names of the researchers” a



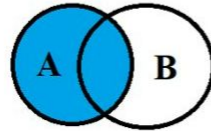
3. Funcionamiento básico

3. Funcionamiento básico

Juntar dos listas a partir de la columna que tienen en común (se añaden el resto de columnas)

Student ID	Name
1001	A
1002	B
1003	C
1004	D

+



Student ID	Department
1004	Mathematics
1005	Mathematics
1006	History
1007	Physics
1008	Computer Science

Student ID	Name	Department
1001	A	NULL
1002	B	NULL
1003	C	NULL
1004	D	Mathematics

3. Funcionamiento básico

Mi lista inicial

Nombre
Sara
Marta
Andrea
Federico
A.

Algoritmo

Nombre	Género
Sara	Female
Sonia	Mostly_female
Andrea	Mostly_female
Federico	Male
Pedro	Male
Andrea	Male



Mi lista final

Nombre	Género
Sara	Female
Marta	Unknown
Andrea	Mostly_female
Federico	Male
A.	Unknown

3. Funcionamiento básico

60% de identificados

El % aceptable de identificados para que tenga sentido los resultados dependen del caso

Muchos algoritmos dan % de probabilidad de que el género haya sido asignado correctamente, o tienen de posibilidades Male/Female/Mostly_male/Mostly_female (o similar)

Mi lista final		
Sara	Female	98%
Marta	Unknown	-
Andrea	Mostly_female	72%
Federico	Male	89%
A.	Unknown	-

4. Cuestiones a considerar

4. Cuestiones a considerar: Países

Mi lista inicial

Nombre	País
Sara	España
Marta	España
Andrea	Italia
Federico	Rusia
A.	Italia

Algoritmo

Nombre	País	Género
Sara	España	Female
Sonia	España	Female
Andrea	España	Mostly_female
Andrea	Italia	Male
Federico	España	Male
Federico	Rusia	Mostly_male
Pedro	España	Male
Alex	Estados Unidos	Unisex
Alex	España	Male

Mi lista final

Nombre	País	Género
Sara	España	Female
Marta	España	Unknown
Andrea	Italia	Male
Federico	Rusia	Mostly_male
A.	Italia	Unknown

4. Cuestiones a considerar: Países

- Nombres con género diferente en países diferentes (Ej. Andrea)
- Nombres eslavos, de Islandia, incluyen información de género en apellidos
- ¿Sobrerrepresentación de ciertos países?

Apellidos eslavos	
Male	Female
-ov	-ova
-ev	-eva
-in	-ina
-ob	-oba
-eb	-eba
-yi	-aya
-yj	-aia
-ky	-iha
-...	-...

4. Cuestiones a considerar: Origen de los datos

¿Sabemos de dónde vienen las listas de nombres que usamos?

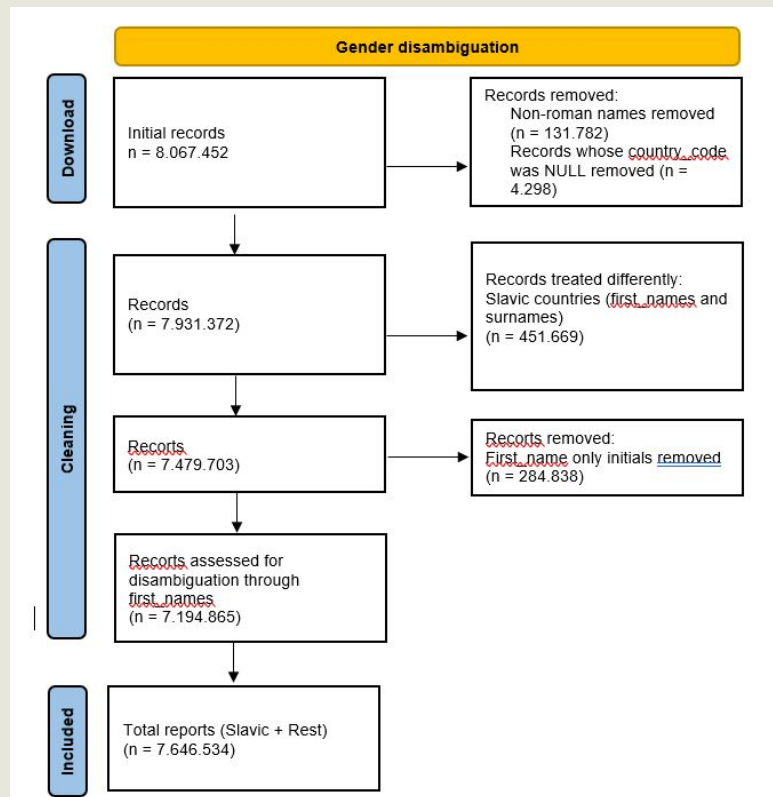
- Nombres occidentales, problemas con Asia principalmente
- Sesgos de bases de datos (i.e., Wikipedia)
- ¿Datos (des)actualizados? (i.e. Nombre "Taylor" en EEUU)

4. Cuestiones a considerar: Datos limpios

Importante tener los datos limpios.

Problemas con:

- Iniciales
- Alfabeto no latino
- Segundos nombres
- Tildes y símbolos raros
- Considerar si queremos desambiguar todos los nombres por igual o por bloques



5. Opciones

¡No dejarse impresionar!
Muchos están repetidos

5. Opciones



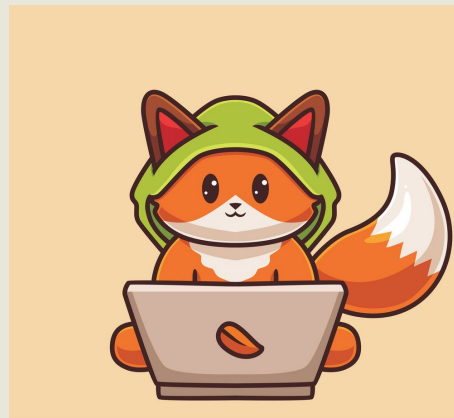
	Gratuitos	Fuentes fiables	Plataforma	Actualizado	Nº nombres
Gender API	✗*	✗	Propia	Sí	6,084,389 de 190 países
Genderize.io	✗	★	Propia	Sí	4.079.646 de 188 países
Gender Guesser	✗	✗	Propia	Sí	9.000.000
Gender-guesser	✓	✗**	Python	No (2016)	44.568
NameAPI	✗*	✗	Propia	?	590.000 de 55 países
Gender-detector	✓	★	Python	No (2015)	125.000 de 2 países
SexMachine	✓	✗**	Python	No (2013)	44,568
World Gender Name Dictionary	✓	★★★	Python	Sí	25.000.000 de 195 territorios

*Free trial

**[Nam_dict.txt](#)

5. Opciones

- Listas de nombres que nosotrxs *matcheamos* con Access, R, etc.
 - [nam_dict.txt](#)
 - Listas caseras
 - Listas desambiguadas de otra gente
- ChatGPT 3.5 ✘
- Usar más de un método (validación)



6. Casos prácticos

6. Casos prácticos: mis datos

	last_name	first_name	country_code
1	Woolf	Virginia	United Kingdom
2	Venturini	Aurora	Argentina
3	Sáez	Juanjo	Spain
4	Pereza	Elisa	Spain
5	Highsmith	Patricia	United States
6	López	Pedro López	Spain
7	Cavarero	Adriana	Italy
8	Yoshimoto	Banana	Japan
9	Innerarity	Daniel	Spain
10	Kulczycki	Emanuel	Poland
11	Harding	Sandra G.	United States
12	Ditlevsen	Tove	Denmark
13	Winsloe	Christa	Germany
14	Tokarczuk	Olga	Poland
15	Bicecci	Verónica Gerber	Mexico
16	Calderón	Javier	Spain
17	Khodyreva	Anastasia A.	Unknown
18	Vallejo	Irene	Spain
19	Hemingway	Ernest	United States
20	O'Farrell	Maggie	United Kingdom
21	Romero	Marga Sánchez	Spain



6. Caso práctico #1: Gender API



- 1** Suba el archivo de Excel o CSV

Navigate through the files or drop an
archivo aquí

Puedes seleccionar y subir un archivo de Excel o CSV de tu escritorio aquí.
Al subir los archivos estás de acuerdo con nuestros términos y condiciones y con nuestra política de privacidad.
Todos los datos se procesan completamente de acuerdo con el GDPR. Lee más aquí.
- 2** Seleccione los campos
- 3** Revise sus datos
- 4** El archivo será procesado
- 5** Descargue el archivo enriquecido

6. Caso práctico #2: World Gender Name Dictionary

¡Cuidado!

	first_name	country_code	level	gender	F	M
2	TT	SA	NaN	not found	not found	not found
6	I. A.	RU	3.0	M	0.0	1.0
0	F.	BR	NaN	not found	not found	not found
5	S.J.	US	NaN	not found	not found	not found
7	J	ES	3.0	M	0.0	1.0
1	S.M	US	NaN	not found	not found	not found
8	R. V.	RU	3.0	M	0.0	1.0
4	M.Z.	UA	NaN	not found	not found	not found
9	G.	FR	2.0	F	1.0	0.0
3	J.L.	GB	NaN	not found	not found	not found

Pasos:

1. [Descargar](#)
2. Abrir en [Google Colab](#) (o similar)
3. Añadir celda 1, 3, 10 y 11 de [este](#) (lo subiré a YSP)
4. Subir vuestro archivo (previamente limpiado) a la carpeta "example"

7. Bibliografía

7. Bibliografía: Cuestiones éticas

- Frohard-Dourlent, H., Dobson, S., Clark, B. A., Doull, M., & Saewyc, E. M. (2017). "I would have preferred more options": Accounting for non-binary youth in health research. *Nursing Inquiry*, 24(1), e12150. <https://doi.org/10.1111/nin.12150>
- González-Salmón, E & Robinson-García, N. (2023). A call for transparency in gender assignment approaches. <https://zenodo.org/records/10036331>
- Heidari, S., Babor, T. F., De Castro, P., Tort, S., & Curno, M. (2016). Sex and gender equity in research: rationale for the SAGER guidelines and recommended use. *Research integrity and peer review*, 1(1), 1-9.
- LaBrada, E. (2016). Categories We Die For: Ameliorating Gender in Analytic Feminist Philosophy. *Publications of the Modern Language Association of America*, 131(2), 449–459. <https://doi.org/10.1632/pmla.2016.131.2.449>
- Lindqvist, A., Sendén, M. G., & Renström, E. A. (2021). What is gender, anyway: A review of the options for operationalising gender. *Psychology & Sexuality*, 12(4), 332–344. <https://doi.org/10.1080/19419899.2020.1729844>
- Medeiros, M., Forest, B., & Öhberg, P. (2020). The Case for Non-Binary Gender Questions in Surveys. *PS: Political Science & Politics*, 53(1), 128–135. <https://doi.org/10.1017/S1049096519001203>
- Mihaljević, H., Tullney, M., Santamaria, L., & Steinfeldt, C. (2019). Reflections on Gender Analyses of Bibliographic Corpora. *Frontiers in Big Data*, 2. <https://www.frontiersin.org/articles/10.3389/fdata.2019.00029>
- Rasmussen, K. C., Maier, E., Strauss, B. E., Durbin, M., Riesbeck, L., Wallach, A., Zamloot, V., & Erena, A. (2019). The Nonbinary Fraction: Looking Towards the Future of Gender Equity in Astronomy (arXiv:1907.04893; Version 1). arXiv. <https://doi.org/10.48550/arXiv.1907.04893>
- Stadler, G., Chesaniuk, M., Haering, S., Roseman, J., Straßburger, V. M., Martina, S., Aisha-Nusrat, A., Maisha, A., Kasia, B., Theda, B., Pichit, B., Marc, D., Sally, D. M., Ruth, D., Ilona, E., Marina, F., Paul, G., Denis, G., Ulrike, G., ... Mine, W. (2023). Diversified innovations in the health sciences: Proposal for a Diversity Minimal Item Set (DiMIS). *Sustainable Chemistry and Pharmacy*, 33, 101072. <https://doi.org/10.1016/j.scp.2023.101072>

7. Bibliografía: Cuestiones técnicas

- Bérubé, N., Ghiasi, G., Sainte-Marie, M., & Larivière, V. (2020). Wiki-Gendersort: Automatic gender detection using first names in Wikipedia. *SocArXiv*. <https://doi.org/10.31235/osf.io/ezw7p>
- Blevins, C., & Mullen, L. (2015). Jane, John ... Leslie? A Historical Method for Algorithmic Gender Prediction. *Digital Humanities Quarterly*, 009(3).
- Karimi, F., Wagner, C., Lemmerich, F., Jadidi, M., & Strohmaier, M. (2016). Inferring Gender from Names on the Web: A Comparative Evaluation of Gender Detection Methods. *Proceedings of the 25th International Conference Companion on World Wide Web - WWW '16 Companion*, 53–54. <https://doi.org/10.1145/2872518.2889385>
- Lax-Martinez, G., Raffo, J. D., & Saito, K. (2023). Identifying the gender of PCT inventors. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.4434107>
- Mryglod, O., Nazarovets, S., & Kozmenko, S. (2023). Peculiarities of gender disambiguation and ordering of non-English authors' names for Economic papers beyond core databases. *Journal of Data and Information Science*, 8(1), 72–89. <https://doi.org/10.2478/jdis-2023-0001>
- Santamaría, L., & Mihaljević, H. (2018). Comparison and benchmark of name-to-gender inference services. *PeerJ Computer Science*, 4, e156. <https://doi.org/10.7717/peerj-cs.156>
- Wais, K. (2016). Gender Prediction Methods Based on First Names with genderizeR. *The R Journal*, 8(1), 17. <https://doi.org/10.32614/RJ-2016-002>

7. Bibliografía: Ejemplos de asignación

- Chan, H. F., & Torgler, B. (2020). Gender differences in performance of top cited scientists by field and country. *Scientometrics*, 125(3), 2421–2447. Scopus. <https://doi.org/10.1007/s11192-020-03733-w>
- El-Ouahi, J., & Larivière, V. (2023). On the lack of women researchers in the Middle East and North Africa. *Scientometrics*, 128(8), 4321–4348. Scopus. <https://doi.org/10.1007/s11192-023-04768-5>
- Fell, C. B., & König, C. J. (2016). Is there a gender difference in scientific collaboration? A scientometric examination of co-authorships among industrial–organizational psychologists. *Scientometrics*, 108(1), 113–141. Scopus. <https://doi.org/10.1007/s11192-016-1967-5>
- Larivière, V., Ni, C., Gingras, Y., Cronin, B., & Sugimoto, C. R. (2013). Bibliometrics: Global gender disparities in science. *Nature*, 504(7479), 211–213.
- Ma, Y., Teng, Y., Deng, Z., Liu, L., & Zhang, Y. (2023). Does writing style affect gender differences in the research performance of articles?: An empirical study of BERT-based textual sentiment analysis. *Scientometrics*, 128(4), 2105–2143. Scopus. <https://doi.org/10.1007/s11192-023-04666-w>

¡Muchas gracias! 🤠

¿Alguna pregunta?



elviragonzalez@go.ugr.es



@egonzalezsalmon