

# MODEL PREDICTIU DE VIOLÈNCIA DE GÈNERE

Autora: María Cobo Ollero  
Tutora: Amanda Fernandez Fontelo



Universitat Autònoma de Barcelona

Grau d'Estadística Aplicada  
Universitat Autònoma de Barcelona

5 de setembre de 2024

# Continguts

- 1 Introducció
- 2 Objectiu
- 3 Dades i metodologia
- 4 Resultats
- 5 Conclusions
- 6 Referències

# Introducció

- ▶ Segons l'ONU, la violència de gènere es refereix als actes nocius dirigits contra una persona o un grup de persones a causa del seu gènere.
- ▶ S'estima que, a escala mundial, aproximadament 1 de cada 3 dones ha estat sotmesa a violència física i/o sexual.

# Objectiu

Aquest treball forma part del projecte *"Data Science Against Gender-Based Violence"* finançat per l'Obra Social La Caixa.

- ▶ L'objectiu d'aquest treball és crear una eina complementària que permeti predir si una dona és una víctima potencial de violència de gènere i que aportï coneixement sobre les característiques de les dones que afavoreixen aquesta situació.

Les dades utilitzades estan extretes de la *Macroencuesta de Violencia contra la Mujer 2019* (MVM) d'Espanya.

- ▶ Única estadística oficial per mesurar la prevalença de violència de gènere a Espanya, són dades públiques
- ▶ Es realitza cada 4 anys, aquesta en 2019
- ▶ A una mostra de 9568 dones residents a Espanya de 16 anys o més
- ▶ Recull les característiques tant de les víctimes com dels agressors
- ▶ El 25% de les dones reconeixen haver patit violència física i/o sexual

S'han aplicat els següents models predictius simples:

- ▶ Naive Bayes (NB)
- ▶ Random Forest (RF)
- ▶ Support Vector Machines (SVM)
- ▶ Simplest Neural Network (SNN)
- ▶ Gradient Boosting (GB)

Per trobar els hiperparàmetres òptims en els diferents models s'ha utilitzat la tècnica *Nested Cross-Validation* (NCV).

Una vegada aplicats els diferents models simples, s'ha aplicat la tècnica *stacked ensemble*. Aquesta combina les prediccions de diversos models d'aprenentatge per a obtenir una predicció final amb un millor rendiment.

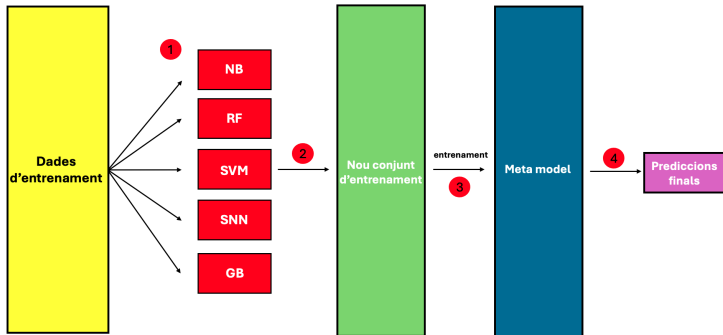


Figura 1: Procediment de l'*stacked ensemble*

Però, s'han seguit diverses línies d'investigació perquè les classes es troben desbalancejades.

- ▶ **Dades originals:** No s'ha tractat el problema del desbalanceig, ens fixem en les mètriques d'interès;
  - Especificitat: Capacitat del model a predir els casos negatius (**violència = si**)
  - Sensibilitat: Capacitat del model a predir els casos positius (**violència = no**)
  - Accuracy equilibrada: Proporció d'observacions correctament classificades que té en compte el desbalanceig de les classes.



- **Dades tractades:** Es tracten les dades per a eliminar el desbalanceig, es fa un *Random under-sampling* (RU) i un SMOTE.

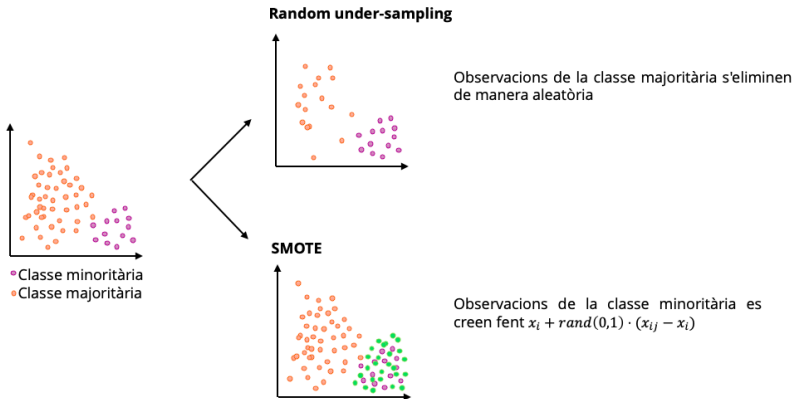


Figura 2: Tractament del desbalanceig de les classes

## Variables importants

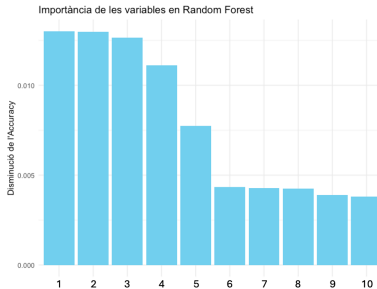


Figura 3: 10 variables més importants

1. Edat
2. Nombre de parelles
3. Edat parella actual
4. Temps relació actual
5. Nombre símptomes últim any
6. Pensament suïcidi últim any
7. Nivell més alt d'estudis oficials
8. Pensament suïcidi
9. Intent suïcidi
10. Situació legal amb la parella

- Factor protector
- Factor de risc

## Variables importants

### Edat:

- ▶ Les dones d'entre 16 i 29 anys (16%) el gairebé 33% reporten haver patit violència.
- ▶ D'entre les dones de 65 anys i més (25%) el quasi 16% reporten haver patit violència.
- ▶ L'edat mitjana, entre les dones que han reportat haver patit violència, és d'aproximadament 46 anys.
- ▶ L'edat mitjana, entre les que no han reportat patir violència, és de quasi 52 anys.

## Variables importants

Nombre de parelles:

- ▶ De les dones que han tingut únicament 1 parella al llarg de la seva vida (56%), el 16% reporten haver patit violència.
- ▶ De les dones que han tingut més de 3 parelles al llarg de la seva vida (7%), el 49% reporten haver patit violència.

## Variables importants

### Edat parella actual:

- ▶ Les dones, les quals l'edat de la seva parella està entre els 16 i 29 anys (38%), el 31% reporten haver patit violència.
- ▶ Les dones, les quals l'edat de la seva parella és de 65 anys i més (16%), el 14% reporten haver patit violència.
- ▶ L'edat mitjana de les parelles actuals, entre les dones que han reportat haver patit violència, és d'aproximadament 47 anys.
- ▶ L'edat mitjana de les parelles actuals, entre les dones que no han reportat patir violència, és de quasi 53 anys.

## Variables importants

Temps de relació amb la parella actual:

- ▶ D'entre les dones que porten menys de 6 mesos amb la seva parella actual (32%), el 50% reporten haver patit violència.
- ▶ D'entre les dones que porten entre 7 i 10 anys de relació (5%), el 33% reporten haver patit violència.
- ▶ D'entre les dones que porten més de 30 anys amb la seva parella (10%), el 14% reporten haver patit violència.

## Variables importants

Nombre símptomes durant l'últim any:

- ▶ Les persones que han patit únicament un símptoma (22%), el 24% reporten haver patit violència.
- ▶ Les persones que han patit 4 símptomes (5%), el 41% reporten haver patit violència.
- ▶ Les persones que han patit tots els símptomes (4%), el 49% reporten haver patit violència.

## Stacked ensemble amb dades originals

Sense aplicar la tècnica *stacked ensemble* s'obté en mitjana; una especificitat del **25%**, una sensibilitat del **94%** i una accuracy equilibrada del **60%**.



MODEL	SPLIT	Especificitat	Sensibilitat	Accuracy equilibrada
NAIVE BAYES	1	48,58	84,10	67,18
	2	50,74	84,19	68,66
	3	49,18	83,04	67,39
	4	46,52	82,17	65,47
	5	47,20	82,31	67,05
	TOTAL	<b>48,44</b>	<b>83,16</b>	<b>67,15</b>
RANDOM FOREST	1	32,03	81,31	62,65
	2	33,26	80,97	63,50
	3	31,22	79,81	62,38
	4	28,69	79,25	60,98
	5	26,60	78,49	60,68
	TOTAL	<b>30,36</b>	<b>79,96</b>	<b>62,04</b>

**Taula 2:** Resultats de fer stacked ensemble amb NB i RF a partir dels models simples amb dades originals

- ▶ Aplicant l'*stacked ensemble* l'especificitat i l'accuracy equilibrada milloren, mentre que la sensibilitat empitjora.
- ▶ S'obtenen millors resultats amb el NB.



## Stacked ensemble amb dades tractades

Abans d'aplicar l'*stacked ensemble* s'observa en mitjana una especificitat del **62%**, una sensibilitat del **66%** i una accuracy del **67%**.



MODEL	SPLIT	Especificitat	Sensibilitat	Accuracy
NAIVE BAYES	1	72,02	71,38	69,53
	2	65,97	71,44	68,60
	3	63,83	72,74	69,74
	4	64,15	73,58	69,53
	5	69,58	71,11	71,58
	TOTAL	<b>67,11</b>	<b>72,05</b>	<b>69,70</b>
RANDOM FOREST	1	72,92	72,25	72,58
	2	70,08	71,65	70,89
	3	65,62	74,21	69,96
	4	65,02	75,48	70,03
	5	70,90	72,31	71,61
	TOTAL	<b>68,91</b>	<b>73,18</b>	<b>71,04</b>

**Taula 3:** Resultats de fer stacked ensemble amb NB i RF a partir dels models simples fent RU.

- ▶ No s'observa cap millora en el rendiment dels models amb RU després d'aplicar la tècnica *stacked ensemble*. Hi ha consistència en els resultats del NB i RF.
- ▶ Continuem explorant perquè en aquest cas aplicar l'*stacked ensemble* no produeix millors prediccions.

## Limitacions

Les limitacions trobades en aquest treball han estat les següents:

- ▶ Desbalanceig en les dades → Seguir explorant altres mètodes, com pot ser, el *Random over-sampling*.
- ▶ Falta informació rellevant sobre característiques de la salut de la dona, com el nombre d'embarassos previs, avortaments, etc.

## Continuació del projecte

Els següents passos en aquest projecte són:

- 1 Repetir la mateixa modelització però estratificant per edats.
- 2 Demanar la MVM a escala europea i reproduir aquesta anàlisi, i respondre a algunes preguntes com són: Hi ha consistència en els resultats? Existeix diferències entre els predictors que prediuen violència en els diferents països europeus? Si és així, quins predictors són els que canvien?
- 3 Quan es publiquin els resultats de la MVM 2024, es reproduirà l'anàlisi.
- 4 Buscar altres tècniques d'*ensemble* per millorar les prediccions.
- 5 Augmentar el nombre de models de classificació simple.

## Conclusions

- ▶ Gràcies a l'eina desenvolupada, ha estat possible augmentar el percentatge de casos de violència predits.
- ▶ Els models de classificació proporcionen idees valuoses per als responsables de polítiques i professionals, ajudant a identificar factors de risc.
- ▶ Cal unificar forces en la lluita contra la violència de gènere i crear una societat més segura i justa per a tothom.

## Referències

- 1 World Health Organization (WHO). (n.d.). *Violence against women prevalence estimates, 2018: Global, regional and national prevalence estimates for intimate partner violence against women and global and regional prevalence estimates for non-partner sexual violence against women*.
- 2 Hastie T., Tibshirani R., Friedman J. (2009). *The elements of statistical learning: Data mining, inference, and prediction (2nd ed.)*. Springer-Verlag. [enllaç](#)
- 3 James, G., Witten, D., Hastie, T., & Tibshirani, R. (2022). *An introduction to statistical learning: With applications in R*. Springer.
- 4 Barton, M., & Lennox, B. (2022). Model stacking to improve prediction and variable importance robustness for soft sensor development. *Digital Chemical Engineering*, 3, 100034. [enllaç](#)
- 5 Arias, M (2023). *Estimating the risk of suffering gender-based violence* [Treball final de màster, Universitat Politècnica de Catalunya & Universitat de Barcelona]