

Inteligencia artificial avanzada para la ciencia de datos 1

Identificación del problema y sus datos

Grupo: 101

Profesor: Gildardo Sánchez

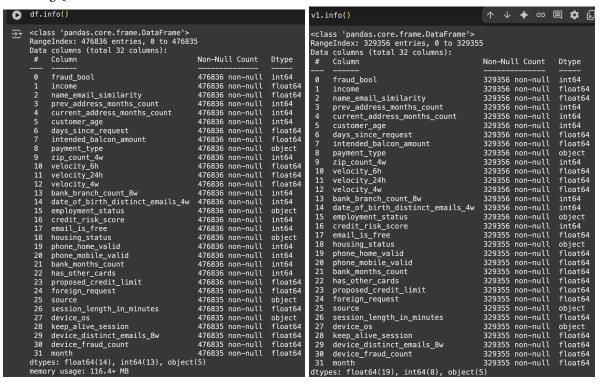
Cristóbal Camarena Hernández - A01642653

agosto 2025

• ¿De que tamaño es el (o los) dataset(s)?



• ¿Qué contiene cada uno?



```
v3.info()
v2.info()
                                                                           <class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.frame.DataFrame'>
                                                                          RangeIndex: 282867 entries, 0 to 282866
Data columns (total 34 columns):
RangeIndex: 278704 entries, 0 to 278703 Data columns (total 34 columns):
                                                                                                                        Non-Null Count
     Column
                                            Non-Null Count
                                                                Dtvpe
                                                                           #
                                                                                Column
                                                                                                                                            Dtype
                                                                                 fraud_bool
                                                                                                                         282867 non-null
                                                                                                                                             int64
      fraud_bool
                                             278704 non-null
                                                                int64
                                                                                                                         282867 non-null
                                                                                                                                             float64
                                                                                 income
     income
                                             278704 non-null
278704 non-null
                                                                float64
     name_email_similarity
                                                                                 name_email_similarity
                                                                                                                        282867 non-null
282867 non-null
                                                                                                                                             float64
                                                                float64
                                                                                 prev_address_months_count
     prev_address_months_count
                                             278704 non-null
                                                                int64
                                                                                                                                             int64
                                                                                 current_address_months_count
                                                                                                                                             int64
                                                                                                                         282867 non-null
     current_address_months_count
                                             278704 non-null
278704 non-null
                                                                int64
                                                                                 customer_age
                                                                                                                         282867 non-null
                                                                                                                                             int64
     customer_age
days_since_request
                                                                int64
                                                                                days_since_request intended_balcon_amount
                                                                                                                        282867 non-null
282867 non-null
                                             278704 non-null
                                                                                                                                             float64
                                                                float64
      intended_balcon_amount
                                             278704 non-null
                                                                float64
                                                                                                                                            float64
                                                                            8
                                                                                 payment_type
                                                                                                                         282867 non-null
                                             278704 non-null
278704 non-null
                                                                                                                                            obiect
     payment_type
                                                                object
                                                                                                                                             int64
                                                                                                                         282867 non-null
     zip_count_4w
velocity_6h
                                                                int64
                                             278704 non-null
                                                                            10
                                                                                 velocity_6h
                                                                                                                         282867 non-null
                                                                                                                                             float64
                                                                float64
     velocity_24h
                                             278704 non-null
                                                                float64
                                                                                 velocity_24h
                                                                                                                        282867 non-null
                                                                                                                                            float64
                                                                                                                         282867 non-null
                                                                                 velocity 4w
                                                                                                                                            float64
     velocity_4w
                                             278703 non-null
                                                                float64
                                                                                                                         282867 non-null
     bank branch count 8w
                                                                                 bank_branch_count_8w
                                                                                                                                             int64
                                             278703 non-null
                                                                float64
                                                                float64
     date_of_birth_distinct_emails_4w
                                            278703 non-null
                                                                            14
                                                                                 date_of_birth_distinct_emails_4w
                                                                                                                        282867 non-null
                                                                                                                                             int64
                                             278703 non-null
                                                                                employment_status
credit_risk_score
     employment_status
                                                                object
                                                                                                                        282867 non-null
                                                                                                                                            object
                                                                                                                        282867 non-null
 16
     credit_risk_score
                                             278703 non-null
                                                                float64
                                                                            16
                                                                                                                                            int64
                                                                                 email_is_free
                                                                                                                         282867 non-null
                                                                float64
                                                                                                                                             int64
     email is free
                                             278703 non-null
                                             278703 non-null
                                                                            18
                                                                                 housing_status
                                                                                                                         282867 non-null
     housing_status
                                                                object
                                                                                                                                            object
                                                                                phone_home_valid phone mobile valid
     phone_home_valid
                                             278703 non-null
                                                                float64
                                                                            19
                                                                                                                        282867 non-null
                                                                                                                                            int64
                                             278703 non-null
278703 non-null
278703 non-null
                                                                float64
float64
                                                                                                                        282867 non-null
 20
     phone_mobile_valid
                                                                            20
                                                                                                                                            int64
                                                                                                                         282867 non-null
                                                                                 bank_months_count
                                                                                                                                             int64
     bank months count
     has_other_cards
                                                                                                                                             int64
                                                                float64
                                                                                 has_other_cards
                                                                                                                         282867 non-null
     proposed_credit_limit
                                                                            23
24
                                                                                 proposed_credit_limit
                                             278703 non-null
                                                                float64
                                                                                                                         282867 non-null
                                                                                                                                             float64
                                                                float64
                                                                                                                         282867 non-null
 24
      foreign_request
                                             278703 non-null
                                                                                 foreign request
                                                                                                                                            int64
                                                                                                                         282867 non-null
 25
                                             278703 non-null
                                                                object
float64
                                                                                 source
                                                                                                                                            object
     source
     session_length_in_minutes
                                             278703 non-null
                                                                                 session_length_in_minutes
                                                                                                                         282867 non-null
 26
                                                                                                                                             float64
     device_os
                                             278703 non-null
                                                                object
                                                                                 device_os
                                                                                                                         282867 non-null
                                                                                                                                            object
                                             278703 non-null
278703 non-null
                                                                            28
29
                                                                                                                         282867 non-null
 28
     keep_alive_session
                                                                float64
                                                                                 keep_alive_session
                                                                                                                                            int64
                                                                float64
                                                                                device_distinct_emails_8w
                                                                                                                         282867 non-null
     device distinct emails 8w
                                                                                                                                             int64
 29
     device_fraud_count
                                             278703 non-null
                                                                float64
                                                                            30
                                                                                 device_fraud_count
                                                                                                                         282867 non-null
                                                                                                                                             int64
     month
                                             278703 non-null
                                                                float64
                                                                                month
                                                                                                                         282867 non-null
                                                                                                                                            int64
                                                                                                                        282867 non-null
282867 non-null
 32
33
     x1
                                             278703 non-null
278703 non-null
                                                                float64
                                                                            32
                                                                                x1
                                                                                                                                            float64
                                                                float64
                                                                            33
                                                                                                                                            float64
     x2
                                                                                x2
 types: float64(24), int64(5), object(5
                                                                                 s: float64(11), int64(18), object
```

```
v4.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 339062 entries, 0 to 339061
Data columns (total 32 columns):
                                         Non-Null Count
#
    Column
                                                           Dtvpe
0
     fraud bool
                                         339062 non-null
                                                           int64
                                         339062 non-null
                                                           float64
     income
     name_email_similarity
                                         339062 non-null
                                                           float64
     prev_address_months_count
                                         339062 non-null
                                                           int64
     current_address_months_count
                                         339062 non-null
                                                           int64
     customer_age
                                         339062 non-null
                                                           int64
     days_since_request intended_balcon_amount
                                         339062 non-null
                                                           float64
                                         339062 non-null
                                                           float64
     payment_type
                                         339062 non-null
                                                           obiect
                                         339062 non-null
     zip_count_4w
                                                           int64
                                         339062 non-null
                                                           float64
 10
     velocity_6h
 11
     velocity_24h
                                         339062 non-null
                                                           float64
 12
     velocity_4w
                                         339062 non-null
                                                           float64
     bank_branch_count_8w date_of_birth_distinct_emails_4w
 13
                                         339062 non-null
                                                           int64
                                         339062 non-null
                                                           int64
 14
     employment_status
                                         339062 non-null
                                                           object
     credit_risk_score
                                         339062 non-null
                                                           int64
                                         339062 non-null
     email_is_free
                                                           int64
     housing_status
                                         339062 non-null
 18
                                                           obiect
     phone_home_valid
 19
                                         339062 non-null
                                                           int64
 20
     phone_mobile_valid
                                         339062 non-null
                                                           int64
 21
     bank_months_count
                                         339062 non-null
                                                           int64
 22
     has other cards
                                         339062 non-null
                                                           int64
     proposed_credit_limit
                                         339062 non-null
 23
                                                           float64
 24
     foreign_request
                                         339062 non-null
                                                           int64
                                         339061 non-null
                                                           object
     session_length_in_minutes
                                         339061 non-null
                                                           float64
     device_os
                                         339061 non-null
                                                           object
 28
     keep_alive_session
                                         339061 non-null
                                                           float64
     device_distinct_emails_8w
                                         339061 non-null float64
 29
     device_fraud_count
                                         339061 non-null
                                                           float64
 30
                                         339061 non-null float64
    month
dtypes: float64(13), int64(14), object(5)
```

• ¿Cuáles son las diferencias entre ellos?

Las principales diferencias son que todos los datasets tienen diferentes filas. También las datasets 2 y 3 tienen más variantes (x1 y x2). Los tipos de datos también varian según el dataset.

32	x1				282867	non-null	float64
33	x2				282867	non-null	float64

• ¿Qué significa el concepto de Fair ML o Fair AI y porque en este contexto resulta muy relevante?

Fair ML / Fair AI es hacer modelos que no traten peor a unas personas que a otras por pertenecer a cierto grupo (edad, ingreso, empleo, etc.). La idea es que la decisión del modelo no dependa injustamente de esos atributos. Por eso hoy evaluar la fairness ya es práctica común en ML, sobre todo en problemas con impacto real en la vida de la gente.

Es importante porque un "positivo" del modelo significa negar abrir una cuenta bancaria. Un falso positivo no es solo un error técnico: le puede cerrar la puerta del sistema financiero a alguien. El propio estudio lo plantea como un dominio "punitivo", donde la decisión afecta directamente el acceso a un servicio básico.

• ¿Implicaciones de usar solo uno de los dataset?

Usar un solo dataset puede generar representatividad limitada, pérdida de patrones, sesgo incrementado, menor robustez, conjunto de variables incompleto, y generalización pobre.