

Proyecto Limpieza de Datos

PROPERATI

- Sankowicz, Javier
- Izquierdo, Stephanie
- Alcoleas, Noelia P.

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

➤ Conocimiento de Base

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

- Conocimiento de Base
- Estandarización de categorías

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

- Conocimiento de Base
- Estandarización de categorías
- Resolución de problemas de formato

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

- Conocimiento de Base
- Estandarización de categorías
- Resolución de problemas de formato
- Asignación de formatos adecuados

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

- Conocimiento de Base
- Estandarización de categorías
- Resolución de problemas de formato
- Asignación de formatos adecuados
- Corrección de valores erróneos

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

- Conocimiento de Base
- Estandarización de categorías
- Resolución de problemas de formato
- Asignación de formatos adecuados
- Corrección de valores erróneos
- Completar datos faltantes

Alcance del Proyecto

A partir de la base de datos recibida, se iniciaron tareas de limpieza y organización de datos. Se pretende lograr una base de datos depurada, que sea eficiente para un posterior análisis y predicción

- Conocimiento de Base
- Estandarización de categorías
- Resolución de problemas de formato
- Asignación de formatos adecuados
- Corrección de valores erróneos
- Completar datos faltantes
- Organización correcta del dataset

El Dataset recibido, tiene 25 variables organizadas en columnas y consta de 121.220 registros en filas.

Variables Originales Presentes en el DataSet:

```
In [3]: data = pd.read_csv('properatti.csv', index_col=0)
        data.shape
```

```
Out[3]: (121220, 25)
```

Operación del aviso : sell, rent

Tipo de la propiedad : house, apartment, ph, store

Nombre del lugar

Nombre del lugar + nombre de sus 'padres'

País

Estado:

ID de geonames del lugar (si está disponible) Latitud,Longitud

Precio original del aviso

Moneda original del aviso : ARS, USD

Precio del aviso en moneda local : ARS

Precio aproximado en USD:

Superficie en m² : referido a la superficie total

Superficie cubierta en m² :

Precio en USD/m² : Precio aproximado en dólares / Superficie Total en m2

Precio por m²: Precio / superficie cubierta en m2

N° de piso, si corresponde:

Ambientes:

Expensas:

URL en Properati

Descripción

Título

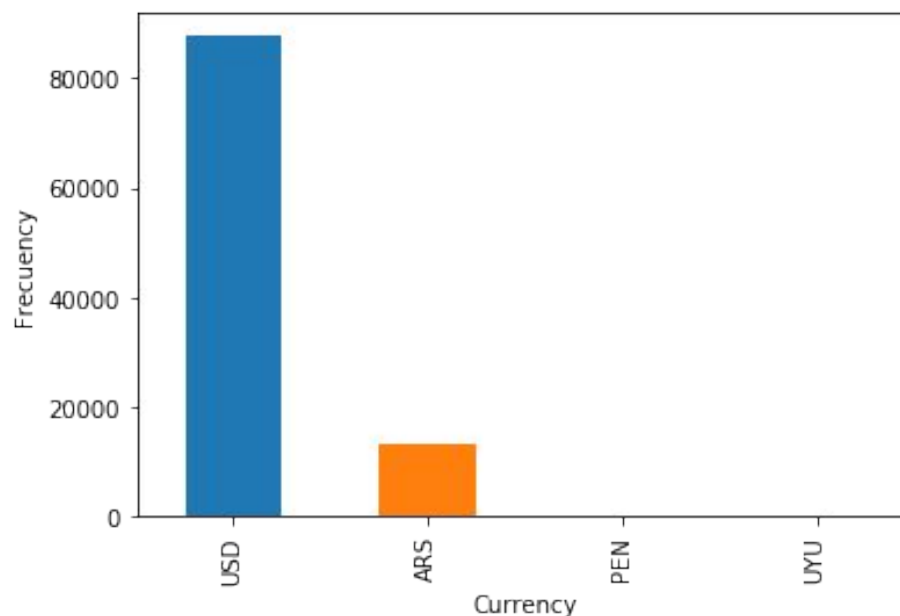
URL de un thumbnail de la primer foto



En la Columna currency, encontramos avisos cotizados en pesos uruguayos y en moneda peruana. Se decidió reclasificar ambos valores a dólares.

```
In [14]: data['currency'].value_counts()
```

```
Out[14]: USD      87587  
ARS      13219  
PEN         2  
UYU         1  
Name: currency, dtype: int64
```



De la revisión de formatos, vimos:

- variables numéricas: datos tipo float
- variables categóricas: datos tipo object

```
In [13]: data.dtypes
```

```
Out[13]: operation          object
property_type             object
place_name                object
place_with_parent_names   object
country_name              object
state_name                object
geonames_id               float64
lat-lon                   object
lat                       float64
lon                       float64
price                     float64
currency                  object
price_aprox_local_currency float64
price_aprox_usd            float64
surface_total_in_m2        float64
surface_covered_in_m2      float64
price_usd_per_m2           float64
price_per_m2              float64
floor                     float64
rooms                     float64
expenses                  float64
properati_url              object
description                object
title                     object
image_thumbnail            object
dtype: object
```



Corrección previa de los datos:

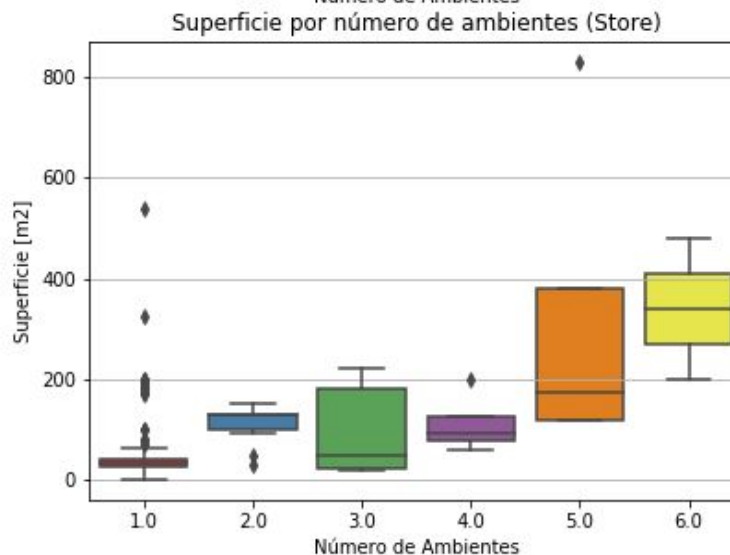
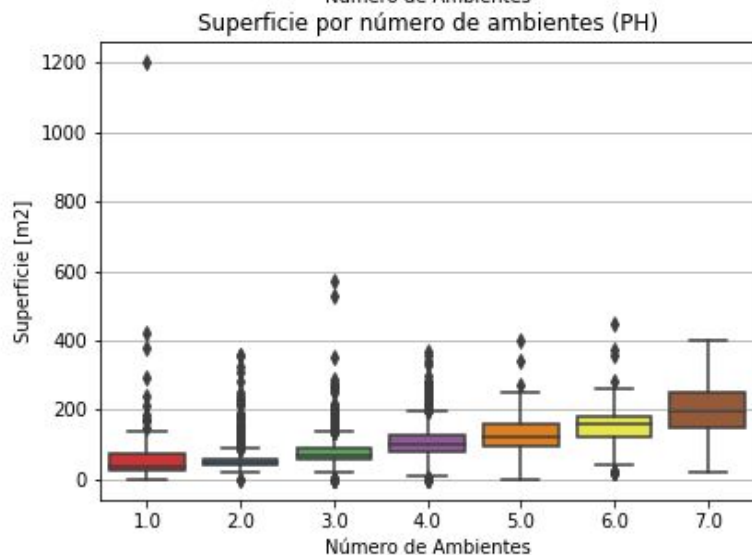
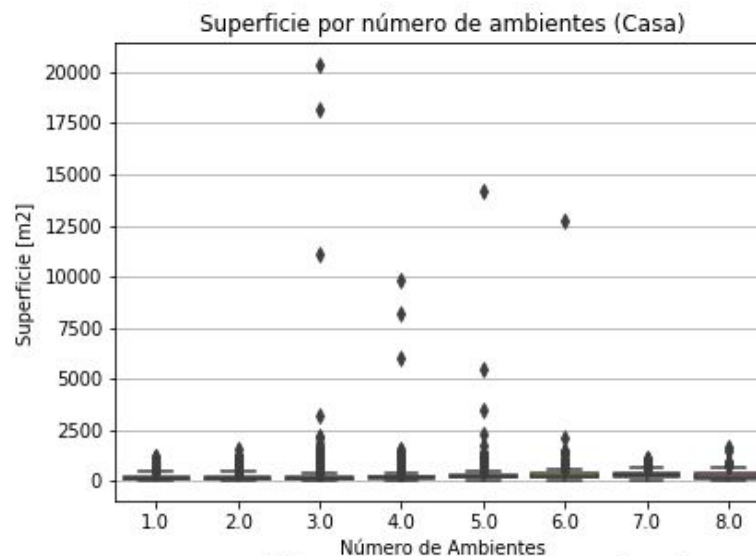
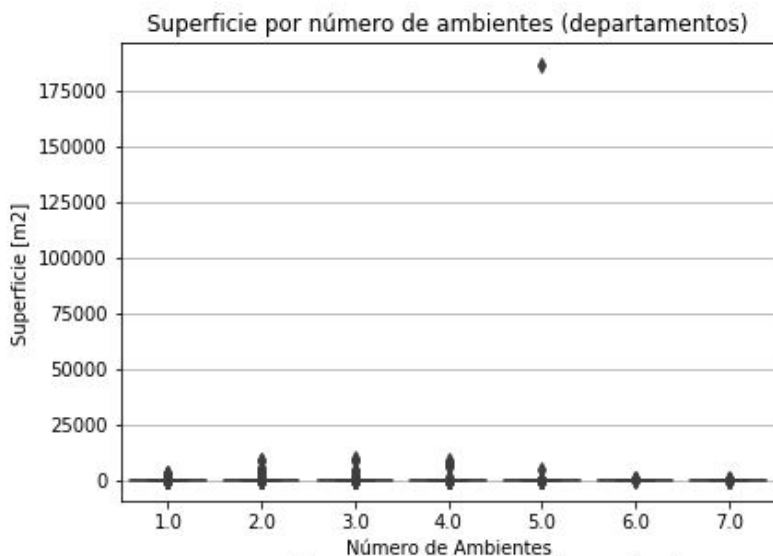
- Surface total in m2: Se determinó una superficie mínima válida.
- Place name: Datos faltantes y datos mal cargados, todos correspondientes a Tigre.
- Currency: Se convirtieron las monedas.
- Validación por Outliers

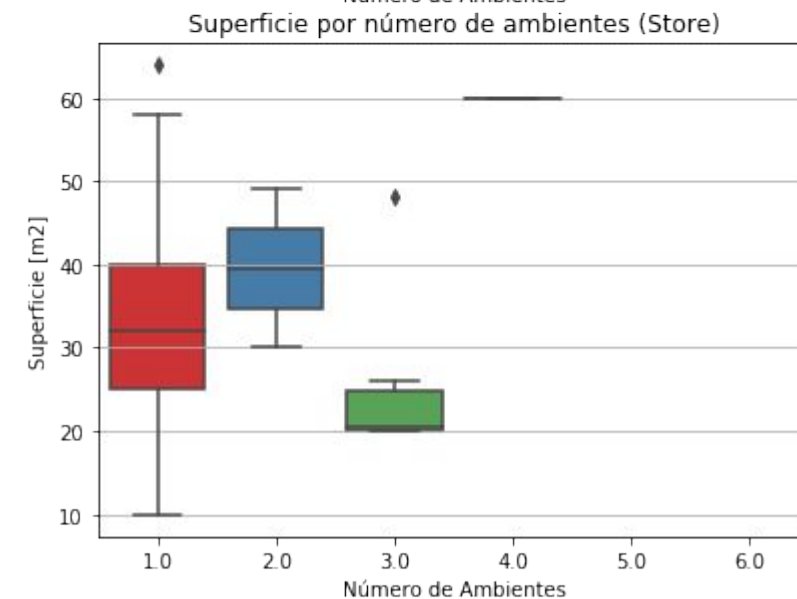
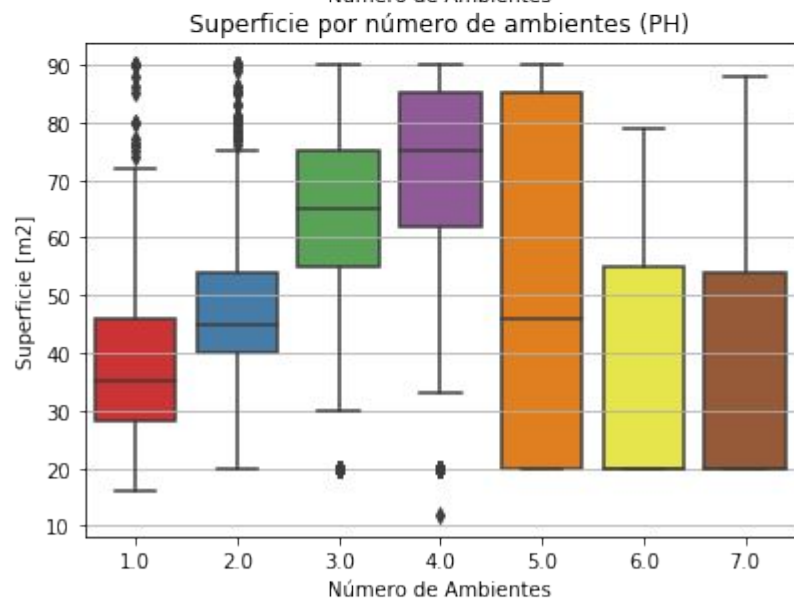
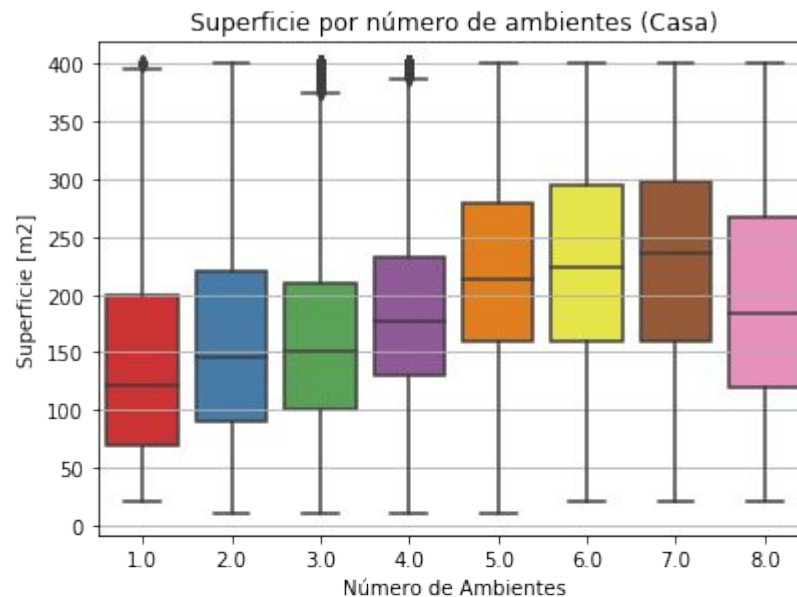
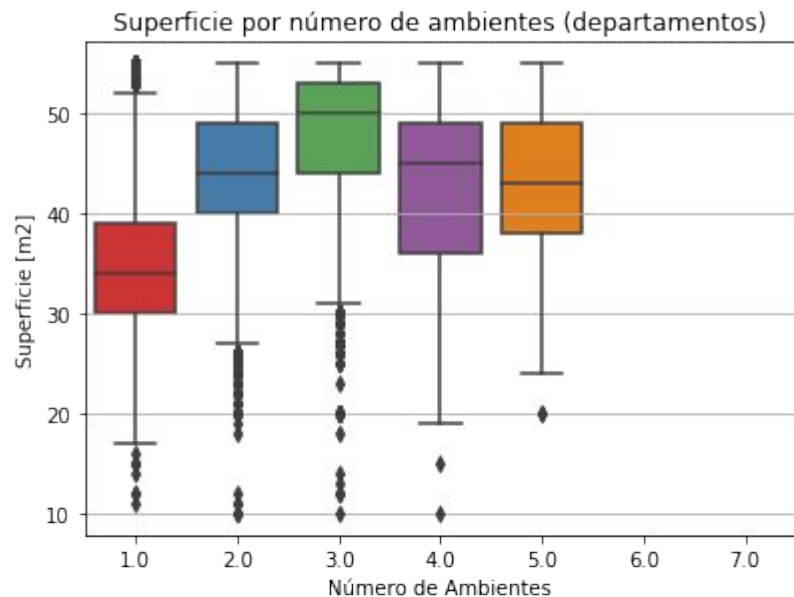


```
data.isnull().sum()
```

operation	0
property_type	0
place_name	23
place_with_parent_names	0
country_name	0
state_name	0
geonames_id	18717
lat-lon	51550
lat	51550
lon	51550
price	20410
currency	20411
price_aprox_local_currency	20410
price_aprox_usd	20410
surface_total_in_m2	39328
surface_covered_in_m2	19907
price_usd_per_m2	52603
price_per_m2	33562
floor	113321
rooms	73830
expenses	106958
properati_url	0
description	2
title	0
image_thumbnail	3112
dtype: int64	

Validación de datos por Outliers Multivariados

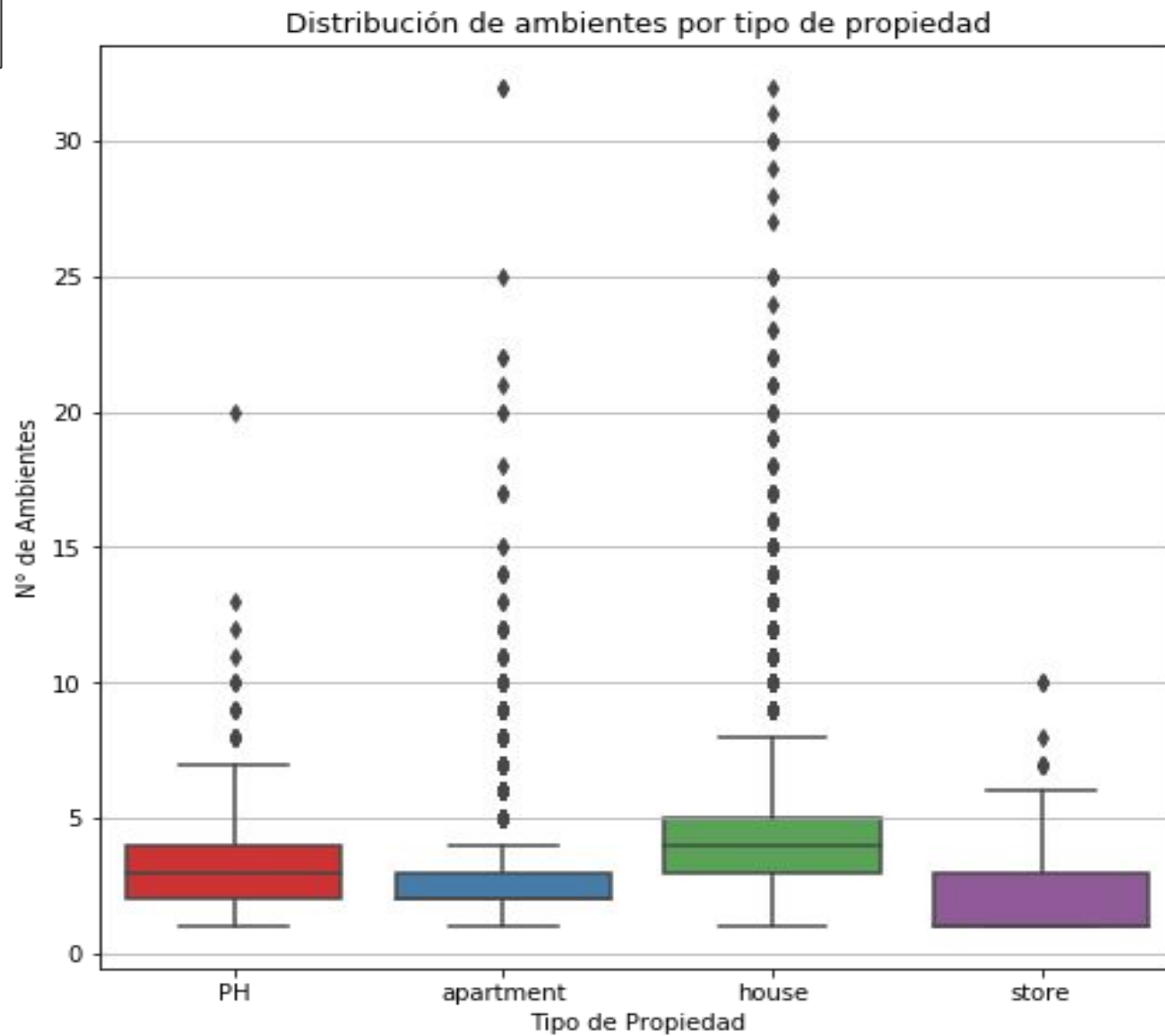




Corrección

Valores Erróneos

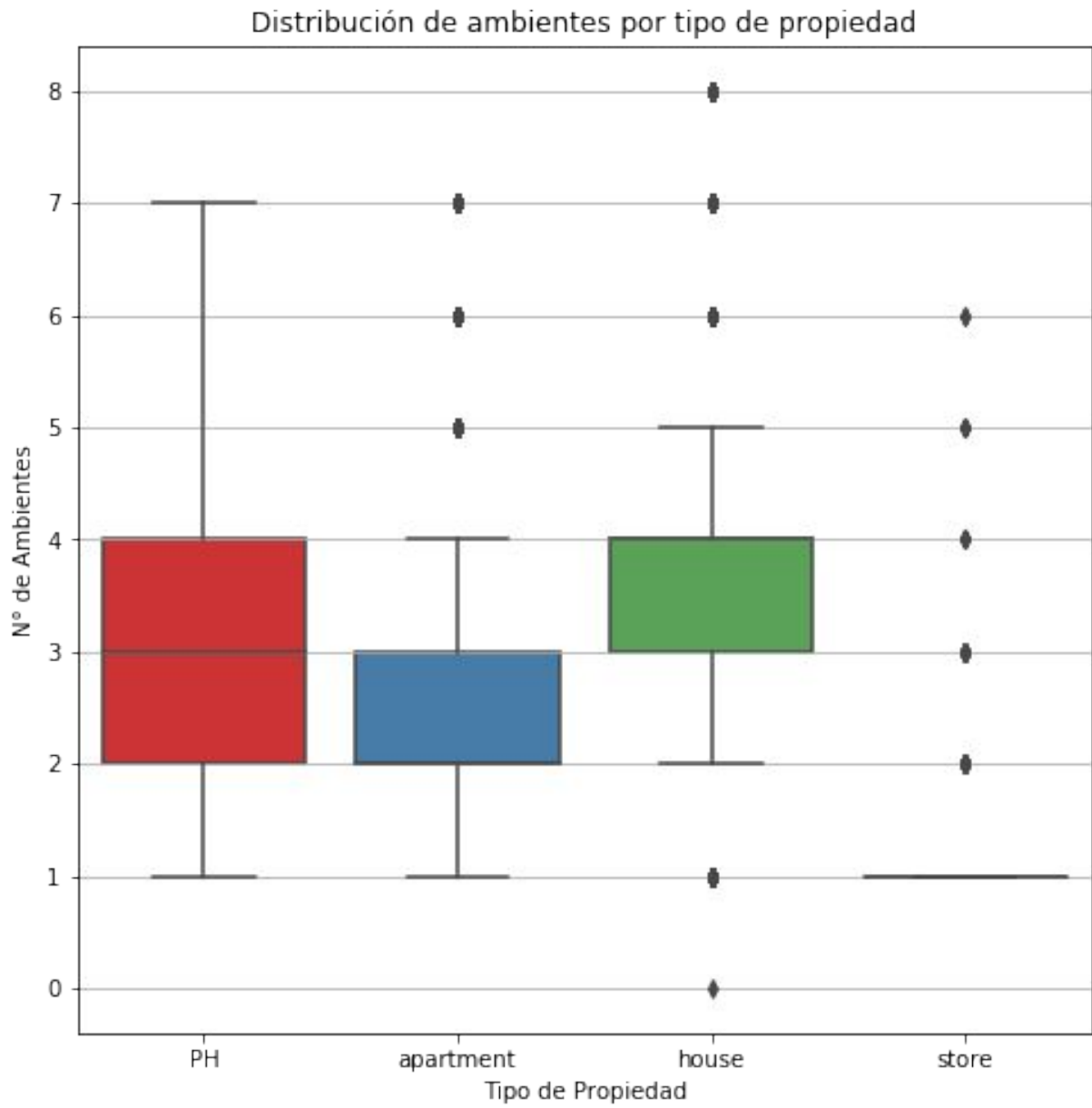
ANTES



Corrección

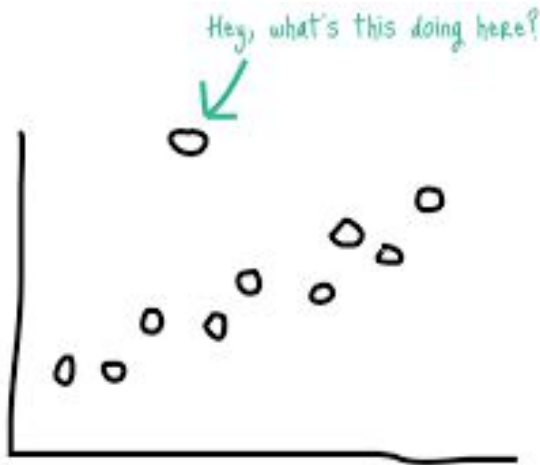
Valores Erróneos

DESPUÉS



Estrategias:

- Búsqueda con Regex
- Álgebra sobre Variables Dependientes
- Índices
- Imputación por media/mediana/moda.

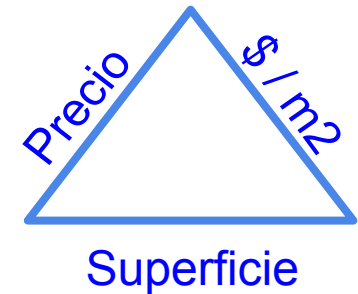


- **Expresiones Regulares:** Búsqueda sobre 'title' y 'description'.



Departamento de 3 Ambientes al Contrafrente. Son 50 M2 totales . Piso 10°. El Edificio cuenta con Grupo Electrógeno. Antigüedad 50 años. Cocina 1.80 X 4.50. Lavadero 1.80 X 1.80, Living-Comedor 3 X 4.50, Dormitorio 2.80 X 3 con placard, Dormitorio 2.50 X 2.60 con placard. Baño Completo. Baulera Independiente en sótano.

- **Álgebra sobre Variables Dependientes:** Aplicadas sobre las distintas variables de precios y superficies presentes en el dataset.



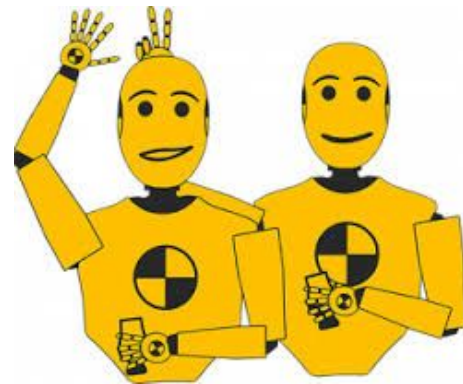
- **Construcción de Índices**

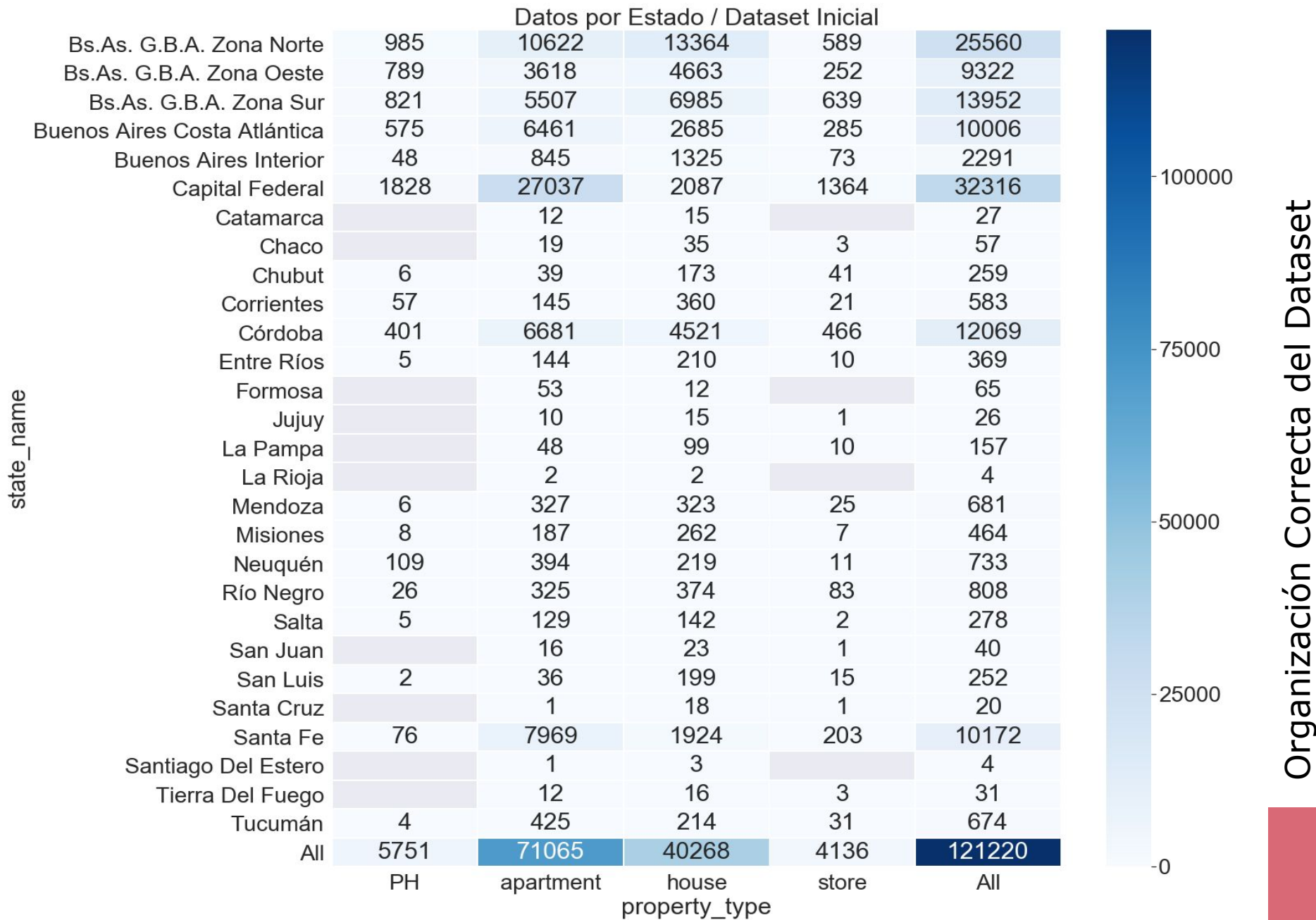
- Índice de construcción por barrio: relaciona las superficies totales y superficies cubiertas por barrio

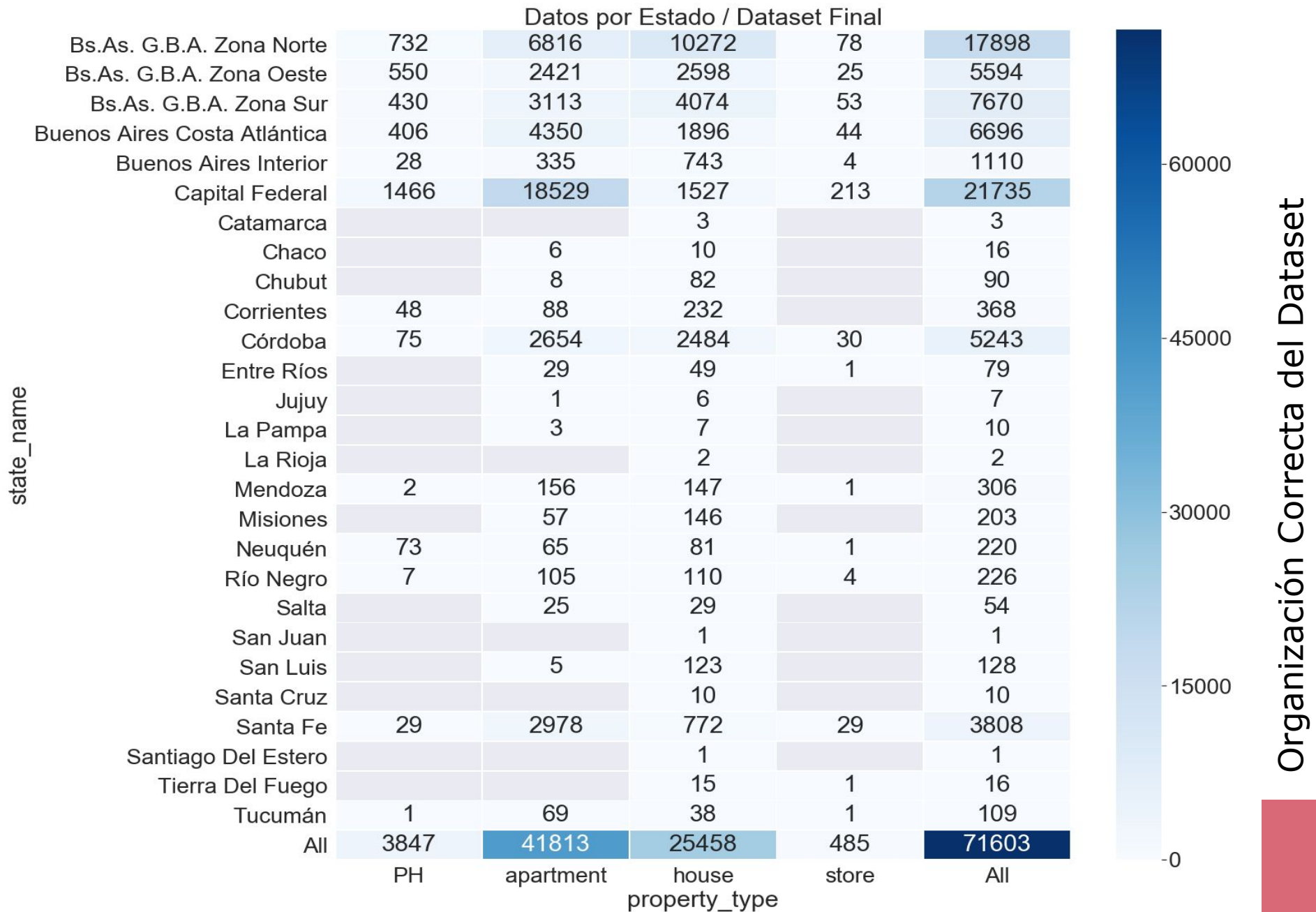
$$\text{Índice de Construcción} = \frac{\text{SuperficieCubierta}}{\text{SuperficieTotal}}$$

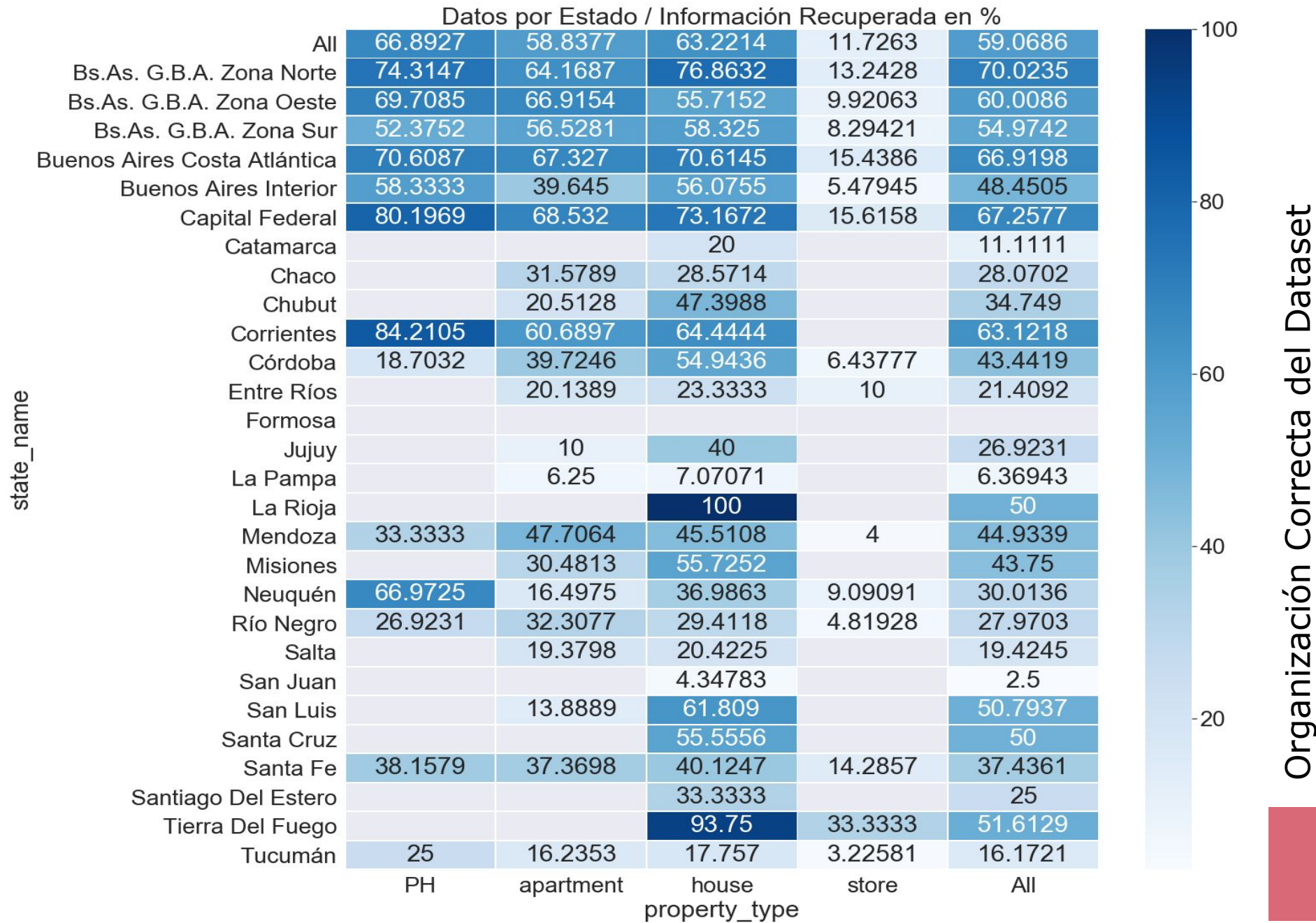
Dataset Final:

- Eliminación de columnas innecesarias.
- Tomamos en cuenta características que aportaran información relevante para la regresión del desafío 2.
- Drop de datos faltantes.
- Variables dummies: Cochera, Pileta, SUM, Patio, Balcón, etc.









¡Gracias!