

File name: VehicleToOrganize.xlsx

| <i>progr</i> | <i>freq</i> | <i>Excel label</i>               | <i>English label</i> | <i>Value</i>                       |
|--------------|-------------|----------------------------------|----------------------|------------------------------------|
| 1            | 100%        | <b>CodVeicolo</b>                | cod_vehicle          | 3                                  |
| 2            | 99%         | <b>anno</b>                      | year                 | 1934                               |
| 3            | 100%        | <b>Marca</b>                     | brand_name           | Alfa Romeo                         |
| 4            | 100%        | <b>Modello_A_Serie</b>           | model_part_a         | 6C 2300                            |
| 5            | 50%         | <b>Modello_B_Versisone</b>       | model_part_b         | Gran Sport                         |
| 6            |             | <b>Mod_Verificato</b>            | verified             |                                    |
| 7            | 100%        | <b>Mod_Famiglia</b>              | family               | 6C                                 |
| 8            | 50%         | <b>Mod_Versione</b>              | version              |                                    |
| 9            | 100%        | <b>Mod_Serie</b>                 | series               | I                                  |
| 10           | 50%         | <b>Mod_Allestimento</b>          | dressing             |                                    |
| 11           | 100%        | <b>Mod_Stage</b>                 | stage                | 1                                  |
| 12           |             | <b>Mod_Potenza</b>               | power                |                                    |
| 13           | 100%        | <b>Modello_Descriz_Originale</b> | description          | 1930 Alfa Romeo 6C 1750 Gran Sport |

1) CodVeicolo: this value is the key\_value. You have to copy this value in new elaborated file without any change.

2) anno: [extract from (13)] (\d{4}) - attention, it's not always first element on the left side; sometimes you can find values like 99 (= 1999), 60' (1960-1969 = we consider 1965) - **definition order = 1**

3) Marca (= brand name): [extract from (13)] - usually is the first text value on the right of [year]. verify existing value into file Brand.xlsx - attention: you can alias values like: Alfa Romeo, Alfa, Alfa Romeo. In any case brand name is omitted and you'll find only family or model name . Verify that [brand\_name] is inside year range [year\_begin ....year\_end inside file brand.xlsx] . **If description contains [clone, replica, handmade, etc...] than, brand\_name = brand-name\_replica**. Copy "Record\_ID" as "ID\_Record\_brand" value into new sheet **definition order = 2**

4) Modello\_A\_Serie: Usually is the same value of [Mod\_Famiglia - Mod\_Serie] (7 - 9). To define [Mod\_Serie] you need to verify inside Model.xlsx file if columns "Mod\_FamigliaA" = Mod\_Famiglia (7). If so, count repetition. If repetition = 1 than Modello\_A\_Serie = [Mod\_FamigliaA] , else find into range [year\_begin ....year\_end inside file Model.xlsx] the value in column Mod\_FamigliaB . Copy "ID\_Record" as "ID\_Record\_model" value into new sheet. **definition order = 4**

5) Modello\_B\_Versisone: all text of point 13, less text of point 2,3,7,4. **definition order = 5**

6) Mod\_Verificato:

7) Mod\_Famiglia: usually is the first text value on the right of [Brand\_name] verify existing value into file Family.xlsx. - Verify that [family] is inside year range [year\_begin ....year\_end inside file Family.xlsx] . Copy "ID\_Record" as "ID\_Record\_family" value into new sheet . **definition order = 3**

8) Mod\_Versione:

9) Mod\_Serie: same than point 4

10) Mod\_Allestimento:

11) Mod\_Stage: column Stage in Model.xlsx of corresponding model value

12) Mod\_Potenza: not to elaborate

13) **Modello\_Descriz\_Originale**: **text to elaborate**

**File name: Brand.xlsx**

| <i>progr</i> | <i>freq</i> | <i>Excel label</i>         | <i>English label</i> | <i>Value</i> |
|--------------|-------------|----------------------------|----------------------|--------------|
| 1            | 100%        | <b>Record_ID</b>           | id                   | 1            |
| 2            | 100%        | <b>Marca</b>               | brand_name           | Lotus        |
| 3            | 10%         | <b>MarcaCompleto</b>       | full_name            |              |
| 4            | 10%         | <b>MarcaCorto</b>          | short_name           |              |
| 5            | 10%         | <b>NomeStoricoCostante</b> | recursive_name       |              |
| 6            | 99%         | <b>DataInizio</b>          | year_begin           | 1952         |
| 7            | 50%         | <b>DataFine</b>            | year_end             |              |

- 1) Record\_ID: this value is the key\_value. You have to copy this value in new elaborated file without any change.
- 2) Marca: always present (=> create a hard code dictionary to add values like: "General Motors", "GM", "General motor corporation", etc.....)
- 3) MarcaCompleto: possible alias
- 4) MarcaCorto: possible alias
- 5) NomeStoricoCostante: possible alias
- 6) DataInizio: if not present, do not evaluate
- 7) DataFine: if not present consider still existing

**File name: Family.xlsx**

| <i>progr</i> | <i>freq</i> | <i>Excel label</i>           | <i>English label</i> | <i>Value</i> |
|--------------|-------------|------------------------------|----------------------|--------------|
| 1            | 100%        | <b>Cod_Famiglia</b>          | family_code          | 39393        |
| 2            | 100%        | <b>ID_Record</b>             | Id                   | 39396        |
| 3            | 100%        | <b>Marca</b>                 | brand_name           | Abarth       |
| 4            | 100%        | <b>Famiglia</b>              | family               | 1500         |
| 5            | 50%         | <b>AnnoFine</b>              | year_begin           |              |
| 6            | 50%         | <b>AnnoInizio</b>            | year_end             |              |
| 7            | 10%         | <b>EsemplariProdotti</b>     | production           |              |
| 8            | 10%         | <b>Fonte_web</b>             | source_web           |              |
| 9            | 100%        | <b>Marca::Record_ID</b>      | id_brand             | 29           |
| 10           | 10%         | <b>Conta_VeicoliPresenti</b> | DB_vehicles          | 1            |
| 11           | 10%         | <b>Conta_Modelli</b>         | count                | 1            |
| 12           | 100%        | <b>AnnoInizioCalcolato</b>   | year_begin_calc      | 1952         |
| 13           | 100%        | <b>AnnoFineCalcolato</b>     | year_end_calc        | 1952         |

- 1) Cod\_Famiglia: this value is the key\_value. You have to copy this value in new elaborated file without any change.
- 2) ID\_Record: this value is the key\_value. You have to copy this value in new elaborated file without any change.
- 3) Marca: brand\_name
- 4) Famiglia: family name
- 5) AnnoFine: year end, if not existing assume value 13 "year\_end\_calculated"
- 6) AnnoInizio: year begin, if not existing assume value 12 "year\_begin\_calculated"
- 7) EsemplariProdotti: do not consider
- 8) Fonte\_web: do not consider

- 9) Marca::Record\_ID: do not consider
- 10) Conta\_VeicoliPresenti: do not consider
- 11) Conta\_Modelli: do not consider
- 12) AnnoInizioCalcolato: calculated year begin
- 13) AnnoFineCalcolato: calculated year begin

**File name: Model.xlsx**

| <i>progr</i> | <i>freq</i> | <i>Excel label</i>         | <i>English label</i> | <i>Value</i> |
|--------------|-------------|----------------------------|----------------------|--------------|
| 1            | 100%        | <b>ID_Record</b>           | id                   | 33791        |
| 2            | 100%        | <b>Marca</b>               | brand_name           | Abarth       |
| 3            | 100%        | <b>Mod_FamigliaA</b>       | model_name_a         | 204          |
| 4            | 100%        | <b>Mod_FamigliaB</b>       | model_name_b         | 204          |
| 5            | 100%        | <b>Serie</b>               | seires               | I            |
| 6            | 100%        | <b>Stage</b>               | stage                | 1            |
| 7            | 100%        | <b>AnnoInizioCalcolato</b> | year_begin           | 1948         |
| 8            | 100%        | <b>AnnoFineCalcolato</b>   | year_end             | 1948         |

- 1) ID\_Record: this value is the key\_value. You have to copy this value in new elaborated file without any change.
- 2) Marca: brand\_name
- 3) Mod\_FamigliaA: = "family"
- 4) Mod\_FamigliaB: = "Mod\_FamigliaA" - "Serie"
- 5) Serie: if count(Mod\_FamigliaA)=1 than "I", else verify "year-begin" and "year\_end" to define "series".
- 6) Stage:
- 7) AnnoInizioCalcolato:
- 8) AnnoFineCalcolato:

**File name: VehcilesOrganized.xlsx**

| <i>progr</i> | <i>freq</i> | <i>Excel label</i>               | <i>English label</i> | <i>Value</i>        |
|--------------|-------------|----------------------------------|----------------------|---------------------|
| 1            | 100%        | <b>CodVeicolo</b>                | id                   | 1                   |
| 2            | 100%        | <b>anno</b>                      | year                 | 1982                |
| 3            | 100%        | <b>Marca</b>                     | brand_name           | Alfa Romeo          |
| 4            | 100%        | <b>Modello_A_Serie</b>           | model_name_a         | F1                  |
| 5            | 50%         | <b>Modello_B_Versisone</b>       | model_name_b         | 182                 |
| 6            |             | <b>Mod_Verificato</b>            | verified             | no                  |
| 7            | 100%        | <b>Mod_Famiglia</b>              | family               | F1                  |
| 8            | 50%         | <b>Mod_Versione</b>              | version              |                     |
| 9            | 100%        | <b>Mod_Serie</b>                 | series               | I                   |
| 10           | 50%         | <b>Mod_Allestimento</b>          | dressing             |                     |
| 11           | 100%        | <b>Mod_Stage</b>                 | stage                | 1                   |
| 12           |             | <b>Mod_Potenza</b>               | power                |                     |
| 13           | 100%        | <b>Modello_Descriz_Originale</b> | description          | 1982 Alfa Romeo 182 |

Same structure of VehicleToOrganize.xlsx

use this file as a reference to the previous conversions of the "description" and the relative distribution of the data in the relevant columns

- 1) CodVeicolo
- 2) anno
- 3) Marca
- 4) Modello\_A\_Serie
- 5) Modello\_B\_Versione
- 6) Mod\_Verificato
- 7) Mod\_Famiglia
- 8) Mod\_Versione
- 9) Mod\_Serie
- 10) Mod\_Allestimento
- 11) Mod\_Stage
- 12) Mod\_Potenza
- 13) Modello\_Descriz\_Originale

File name: **AI\_VehcilesOrganized.xlsx**

| <i>progr</i> | <i>freq</i> | <i>Excel label</i>               | <i>English label</i> | <i>Value</i>        |
|--------------|-------------|----------------------------------|----------------------|---------------------|
| 1            | 100%        | <b>CodVeicolo</b>                | id                   | 1                   |
| 2            | 100%        | <b>anno</b>                      | year                 | 1982                |
| 3            | 100%        | <b>Marca</b>                     | brand_name           | Alfa Romeo          |
| 4            | 100%        | <b>Modello_A_Serie</b>           | model_name_a         | F1                  |
| 5            | 50%         | <b>Modello_B_Versione</b>        | model_name_b         | 182                 |
| 6            |             | <b>Mod_Verificato</b>            | verified             | no                  |
| 7            | 100%        | <b>Mod_Famiglia</b>              | family               | F1                  |
| 8            | 50%         | <b>Mod_Versione</b>              | version              |                     |
| 9            | 100%        | <b>Mod_Serie</b>                 | series               | I                   |
| 10           | 50%         | <b>Mod_Allestimento</b>          | dressing             |                     |
| 11           | 100%        | <b>Mod_Stage</b>                 | stage                | 1                   |
| 12           |             | <b>Mod_Potenza</b>               | power                |                     |
| 13           | 100%        | <b>Modello_Descriz_Originale</b> | description          | 1982 Alfa Romeo 182 |
| 14           | 100%        | <b>Cod_brand</b>                 | Cod_brand            |                     |
| 15           | 100%        | <b>Cod_family</b>                | Cod_family           |                     |
| 16           | 100%        | <b>Cod_model</b>                 | Cod_model            |                     |
| 17           | 100%        | <b>Accuracy</b>                  | Accuracy             |                     |

Accuracy:

- A = 100% (Brand, family, model is inside year-range and in description)
- B = 99% (Brand, family, model is inside year-range **but not** in description)
- C = 80% (Brand, family, model is inside year-range and in description)
- D = 50% (Brand, family, in description but year range **is not inside range** +/- 2 year)
- E = 40% (Brand, family, in description but year range **is not inside range MORE** +/- 2 year)

F = 25% (Only Year and Brand in description, **and** inside year range))  
G = 20% (Only Year and Brand in description, **but not** inside year range)  
H = 5% (Family is not listed in family)  
J = 0% (Brand is not listed in brand.xlsx)