

# Preparazione dei Dati

Il notebook si occupa di trasformare i dati grezzi messi a disposizione dal [task di Kaggle](#): il task di preparazione si propone di:

- fornire una prima analisi generale dei dataset.
- effettuare una pulizia di dati mancanti, ridondanti o poco significativi.
- aggiungere in forma esplicita nuova informazione utile.
- trasformare i dati affinché possano essere processati in maniera corretta da un algoritmo di Machine-Learning.

## Letture dei dati

```
In [1]: # Librerie esterne
import math
import os
import warnings

import pandas as pd
import numpy as np

from sklearn.preprocessing import OneHotEncoder
from sklearn.model_selection import train_test_split

from datetime import datetime as dt

warnings.filterwarnings('ignore')
```

Letture dei dataset forniti da Kaggle.

```
In [2]: # Local file paths
dir_name = 'datasets'

fp_properties2016 = dir_name + "/properties_2016.csv"
fp_properties2017 = dir_name + "/properties_2017.csv"
fp_train2016 = dir_name + "/train_2016.csv"
fp_train2017 = dir_name + "/train_2017.csv"
```

```
In [3]: # Lettura dei dataframe
df_properties2016 = pd.read_csv(fp_properties2016, low_memory=False)
df_train2016 = pd.read_csv(fp_train2016, low_memory=False)
df_properties2017 = pd.read_csv(fp_properties2017, low_memory=False)
df_train2017 = pd.read_csv(fp_train2017, low_memory=False)
```

```
In [4]: # Dimensionalità
print(f'Properties 2016 {df_properties2016.shape}')
print(f'Train 2016 {df_train2016.shape}')
print(f'Properties 2017 {df_properties2017.shape}')
print(f'Train 2017 {df_train2017.shape}')
```

Properties 2016 (2985217, 58)  
Train 2016 (90279, 3)  
Properties 2017 (2985217, 58)  
Train 2017 (77613, 3)

Non dispongono del log-error di ogni casa, ma solo di quelle che sono state vendute: selezione solo l'insieme di case di cui ho a disposizione il log-error.

Unione in un unico dataset: metting le sole case di cui conosco il log-error.

Se una casa ha più log-error, la colonna è copiata e abbinata a ciascuna data di vendita.

```
In [5]: # Right-join
df_2016 = pd.merge(df_properties2016, df_train2016, how='right', left_on='parcelid', right_on='parcelid')
df_2017 = pd.merge(df_properties2017, df_train2017, how='right', left_on='parcelid', right_on='parcelid')
```

```
In [6]: # Dimensionalità
print(f'Properties 2016 {df_2016.shape}')
print(f'Properties 2017 {df_2017.shape}')
```

Properties 2016 (30273, 60)  
Properties 2017 (77613, 60)

Unisco in un unico dataset i dati del 2016 e del 2017.

```
In [7]: # Concat
dfAll = pd.concat([df_2016, df_2017], ignore_index=True)
```

```
In [8]: del(df_2016, df_2017)
```

```
In [9]: dfAll.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 167888 entries, 0 to 167887
Data columns (total 60 columns):
#   Column                                Non-Null Count  Dtype
---  ---
0   parcelid                             167888 non-null  int64
1   airconditioningtypeid                53788 non-null  float64
2   architecturalstyletypeid             468 non-null    float64
3   basementsqft                        93 non-null     float64
4   bathroomcnt                         167854 non-null float64
5   bedroomcnt                          167854 non-null float64
6   buildingclastypeid                  31 non-null     float64
7   buildingqualitytypeid               107173 non-null float64
8   calculatedbathnbr                   166056 non-null float64
9   decktypeid                          12893 non-null  float64
10  finishedfloorisquarefeet            12893 non-null  float64
11  calculatedfinishedsquarefeet        166992 non-null float64
12  finishedsquarefeet12                 75 non-null     float64
13  finishedsquarefeet15                 691 non-null    float64
14  finishedsquarefeet50                 12893 non-null float64
15  finishedsquarefeet150               807 non-null    float64
16  finishedsquarefeet6                  807 non-null    float64
17  fips                                  167854 non-null float64
18  fireplacecnt                        17896 non-null  float64
19  fullbathcnt                         166056 non-null float64
20  garagecarcnt                        55457 non-null  float64
21  garagetotalsqft                     55457 non-null  float64
22  hashottuborspa                      3904 non-null   object
23  heatingandsystemtypeid              105651 non-null float64
24  latitude                             167854 non-null float64
25  longitude                             167854 non-null float64
26  lotsizeisquarefeet                  149446 non-null float64
27  poolcnt                              34075 non-null  float64
28  poolsizeusm                          1838 non-null   float64
29  pooltypeid0                          1626 non-null   float64
30  pooltypeid2                          2278 non-null   float64
31  pooltypeid7                          31776 non-null  float64
32  propertycountylandusecode           167853 non-null object
33  propertylandusetypeid               167854 non-null float64
34  propertyzoningdesc                  108789 non-null object
35  rawconstructionblock                167854 non-null float64
36  regionidcity                        164579 non-null float64
37  regionidcounty                      167854 non-null float64
38  regionidneighborhood                66986 non-null float64
39  regionidzip                          167769 non-null float64
40  roomcnt                             167854 non-null float64
41  storytypeid                          93 non-null     float64
42  threequarterbathnbr                 22115 non-null  float64
43  typeconstructiontypeid              522 non-null    float64
44  unitcnt                             109056 non-null float64
45  yardbuildingsqft17                  503 non-null    float64
46  yardbuildingsqft26                  165 non-null    float64
47  yearbuilt                           166828 non-null float64
48  numberofstories                     38169 non-null  object
49  fireplaceflag                       394 non-null    object
50  structuretaxvaluedollarcnt          167359 non-null float64
51  taxvaluedollarcnt                   167852 non-null float64
52  assessmentyear                      167854 non-null float64
53  landtaxvaluedollarcnt               167851 non-null float64
54  taxamount                           167843 non-null object
55  taxdelinquencyflag                  4683 non-null   object
56  taxdelinquencyyear                  4683 non-null   object
57  censustractblock                    167888 non-null float64
58  logerror                             167888 non-null object
59  transactiondate                     167888 non-null object
dtypes: float64(53), int64(1), object(6)
memory usage: 76.9+ MB
```

## Casting dei tipi

Prima di processare i dati è effettuato un casting tutti i tipi di dato numerici da 64 bit (tipo di memoria) a 32 bit con il doppio scopo di ridurre l'ingombro dei dati e effettuare calcoli più efficienti.

```
In [10]: # Given a dataframe cast all numeric type from 64 bit to 32 bit
def int_float_to32(df):
    for c, dtype in zip(df.columns, df.dtypes):
        if dtype == np.float64:
            df[c] = df[c].astype(np.float32)
        elif dtype == np.int64:
            df[c] = df[c].astype(np.int32)
    return df
```

```
In [11]: dfAll = int_float_to32(dfAll)
```

```
In [12]: dfAll.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 167888 entries, 0 to 167887
Data columns (total 60 columns):
#   Column                                Non-Null Count  Dtype
---  ---
0   parcelid                             167888 non-null  float32
1   airconditioningtypeid                53788 non-null  float32
2   architecturalstyletypeid             468 non-null    float32
3   basementsqft                        93 non-null     float32
4   bathroomcnt                         167854 non-null float32
5   bedroomcnt                          167854 non-null float32
6   buildingclastypeid                  31 non-null     float32
7   buildingqualitytypeid               107173 non-null float32
8   calculatedbathnbr                   166056 non-null float32
9   decktypeid                          12893 non-null  float32
10  finishedfloorisquarefeet            12893 non-null  float32
11  calculatedfinishedsquarefeet        166992 non-null float32
12  finishedsquarefeet12                 75 non-null     float32
13  finishedsquarefeet15                 691 non-null    float32
14  finishedsquarefeet50                 12893 non-null float32
15  finishedsquarefeet150               807 non-null    float32
16  finishedsquarefeet6                  807 non-null    float32
17  fips                                  167854 non-null float32
18  fireplacecnt                        17896 non-null  float32
19  fullbathcnt                         166056 non-null float32
20  garagecarcnt                        55457 non-null  float32
21  garagetotalsqft                     55457 non-null  float32
22  hashottuborspa                      3904 non-null   object
23  heatingandsystemtypeid              105651 non-null float32
24  latitude                             167854 non-null float32
25  longitude                             167854 non-null float32
26  lotsizeisquarefeet                  149446 non-null float32
27  poolcnt                              34075 non-null  float32
28  poolsizeusm                          1838 non-null   float32
29  pooltypeid0                          1626 non-null   float32
30  pooltypeid2                          2278 non-null   float32
31  pooltypeid7                          31776 non-null  float32
32  propertycountylandusecode           167853 non-null object
33  propertylandusetypeid               167854 non-null float32
34  propertyzoningdesc                  108789 non-null object
35  rawconstructionblock                167854 non-null float32
36  regionidcity                        164579 non-null float32
37  regionidcounty                      167854 non-null float32
38  regionidneighborhood                66986 non-null float32
39  regionidzip                          167769 non-null float32
40  roomcnt                             167854 non-null float32
41  storytypeid                          93 non-null     float32
42  threequarterbathnbr                 22115 non-null  float32
43  typeconstructiontypeid              522 non-null    float32
44  unitcnt                             109056 non-null float32
45  yardbuildingsqft17                  503 non-null    float32
46  yardbuildingsqft26                  165 non-null    float32
47  yearbuilt                           166828 non-null float32
48  numberofstories                     38169 non-null  object
49  fireplaceflag                       394 non-null    object
50  structuretaxvaluedollarcnt          167359 non-null float32
51  taxvaluedollarcnt                   167854 non-null float32
52  assessmentyear                      167851 non-null float32
53  landtaxvaluedollarcnt               167843 non-null object
54  taxamount                           167843 non-null object
55  taxdelinquencyflag                  4683 non-null   object
56  taxdelinquencyyear                  4683 non-null   object
57  censustractblock                    167888 non-null float32
58  logerror                             167888 non-null object
59  transactiondate                     167888 non-null object
dtypes: float32(53), int32(1), object(6)
memory usage: 42.3+ MB
```

Il casting è avvenuto in maniera corretta.

```
In [13]: dfAll.shape
```

```
Out[13]: (167888, 60)
```

## Prima analisi delle occorrenze dei parcelid

Analisi delle occorrenze dei parcelid: quante volte una casa è stata venduta tra 2016 e 2017.

```
In [14]: dfAll.loc[:, 'parcelid'].value_counts().head(20)
```

```
Out[14]:
10857130    3
1991059     3
11842707    3
14010551    3
12478591    3
14672826    3
17164212    3
17237150    3
12612211    3
1186136     2
1991474     2
12273662    2
11061551    2
14659784    2
12467034    2
1256654     2
14008322    2
12752161    2
11887100    2
12239653    2
Name: parcelid, dtype: int64
```

Una casa è stata venduta al massimo tre volte: estraggo le case vendute tre volte tra 2016 e 2017.

```
In [15]: houses = list(dfAll.loc[:, 'parcelid'].value_counts().
                    dfAll.loc[:, 'parcelid'].value_counts() == 3
                    ).to_dict().keys()
```

```
Out[15]: {10857130,
1991059,
11842707,
14010551,
12478591,
14672826,
17164212,
17237150,
12612211,
1186136,
1991474,
12273662,
11061551,
14659784,
12467034,
1256654,
14008322,
12752161,
11887100,
12239653}
```

I numeri di riga sono ora mescolati: ripristino del numero di riga.

```
In [16]: # Given a dataframe set its rows in range from 0 to n in ascending order
def arrange_rows(df):
    df.index = np.arange(len(df))
    return df
```

```
In [17]: for df in [X_train, X_val, X_test, y_train, y_val, y_test]:
            df = arrange_rows(df)
```

```
In [18]: X_train.loc[:, ['parcelid', 'transactiondate']].head()
```

```
Out[18]:
   parcelid  transactiondate
0  14217523      2017-08-02
1  11199964      2017-07-11
2  12627031      2016-04-15
3  13992985      2017-05-02
4  12086693      2016-10-13
```

```
In [19]: y_train.head()
```

```
Out[19]:
0    0.057681
1   -0.010615
2    0.020800
3    0.001967
4   -0.020200
Name: logerror, dtype: float32
```

I numeri di riga sono ora mescolati: ripristino del numero di riga.

```
In [20]: # Given a dataframe set its rows in range from 0 to n in ascending order
def arrange_rows(df):
    df.index = np.arange(len(df))
    return df
```

```
In [21]: for df in [X_train, X_val, X_test, y_train, y_val, y_test]:
            df = arrange_rows(df)
```

```
In [22]: X_train.loc[:, ['parcelid', 'transactiondate']].head()
```

```
Out[22]:
   parcelid  transactiondate
0  14217523      2017-08-02
1  11199964      2017-07-11
2  12627031      2016-04-15
3  13992985      2017-05-02
4  12086693      2016-10-13
```

```
In [23]: y_train.head()
```

```
Out[23]:
0    0.057681
1   -0.010615
2    0.020800
3    0.001967
4   -0.020200
Name: logerror, dtype: float32
```

## Rappresentazione non corretta dei Nan

Analizzando il dataset si evince che alcune feature rappresentano la assenza di una caratteristica con un Nan, quando in realtà ai fini di algoritmi di Machine Learning sarebbe più corretta una rappresentazione con valori zero o False, ad esempio:

- fireplaceflag ha valori Nan e True, sarebbe opportuno la conversione in una variabile binaria.
- fireplace e poolcnt: hanno un valore numerico se l'elemento è presente, Nan se non è presente. Conversione della rappresentazione dell'assenza con 0.

```
In [30]: # Given a dataframe and a column name, column's values are set to zero if Nan, one otherwise
def set_zero_one(df, col_name):
    for col_name in col_names:
        is_na = df.loc[:, col_name].isna()
        df.loc[:, col_name][is_na] = 0.
    return df
```

```
In [31]: # Given a dataframe and a column name, values of that column are set to zero if Nan
def nan_to_zero(df, col_name):
    for col_name in col_names:
        df.loc[:, col_name].fillna(0., inplace=True)
    return df
```

```
In [32]: for X in [X_train, X_val, X_test]:
            X = set_zero_one(X, 'fireplaceflag')
            X = nan_to_zero(X, ['fireplacecnt', 'poolcnt'])
```

## Rimozione degli Outlier

Rimuovo dal Train righe che presentano logerror estremi: potrebbero costituire dei punti di rumore e inficiare un corretto funzionamento degli algoritmi di Machine-Learning.

```
In [33]: # Given a dataframe and a column remove all rows which target value is under the first or over the last percentile
def remove_outlier(X, y):
    for col_name in col_names:
        q1 = y < np.percentile(y, 99.5)
        q3 = y > np.percentile(y, 0.01)
        out = list(map(lambda col, q1, q2: out1 and out2, out1, out2))
        X = X[out]
        y = y[out]
    return X, y
```

```
In [34]: dimensionality()
```

```
Out[34]: X_train (99724, 33)
X_val (33578, 33)
X_test (33578, 33)
```

```
In [35]: X_train, y_train = remove_outlier(X_train, y_train)
```

```
Out[35]: dimensionality()
```

```
Out[36]: X_train (99724, 33)
X_val (33578, 33)
X_test (33578, 33)
y_val (33578, 1)
y_test (33578, 1)
```

Sono state rimosse circa un migliaio di righe dal Train.

## Rimozione delle colonne con alta percentuale di Nan

Rimozione delle colonne con un'alta percentuale di valori assenti: queste arricchiscono l'informazione del dataset in maniera molto limitata.

```
In [37]: # Given the dataframe and the name of a column returns the column
def get_col(df, col_name):
    return df.loc[:, col_name]
```

```
In [38]: # Given a dataframe remove columns with Nan-count and Nan-percentage
def get_col_nan_info(col):
    count = col.isna().sum()
    tot = len(col)
    perc = count/tot
    return count, perc
```

```
In [39]: # Given the df and a cut-off returns a list of column names with Nan-percentage greater or equal to the cut-off
def get_cols_over_nan_percentage(df, a_cut_off):
    names = df.columns
    overPercentage = []
    for name in names:
        col = get_col(df, name)
        _, perc = get_col_nan_info(col)
        if perc > a_cut_off:
            overPercentage.append(name)
    return overPercentage
```

```
In [40]: col_to_delete = get_cols_over_nan_percentage(X_train, 0.6)
```

```
Out[40]: for o in col_to_delete:
            print(f'col: {get_col_nan_info(get_col(X_train, o))}')
            print(f'Length: {len(col_to_delete)}')
```

```
Out[41]: airconditioningtypeid: (67655, 0.678422445908242)
architecturalstyletypeid: (99455, 0.997302550519433)
basementsqft: (99676, 0.999186715334233)
buildingclastypeid: (49709, 0.998939792643178108)
decktypeid: (99016, 0.992900405118126)
finishedfloorisquarefeet: (91940, 0.921944570049337)
finishedsquarefeet12: (98676, 0.9895186715346323)
finishedsquarefeet15: (95882, 0.961473673218082)
finishedsquarefeet50: (91940, 0.921944570049337)
finishedsquarefeet150: (99264, 0.998268620939)
garagecarcnt: (66611, 0.6679535518029762)
garagetotalsqft: (66611, 0.6679535518029762)
hashottuborspa: (97994, 0.976655140186961)
poolsizeusm: (98619, 0.9889194175925554)
pooltypeid0: (98711, 0.988941638201436)
pooltypeid2: (98407, 0.9867935018984849)
pooltypeid7: (80817, 0.810406722542497)
regionidneighborhood: (59904, 0.600979262765232)
storytypeid: (98676, 0.9895186715346323)
threequarterbathnbr: (86493, 0.867328137258934)
typeconstructiontypeid: (99417, 0.9968215033492439)
yardbuildingsqft17: (86465, 0.921944570049337)
yardbuildingsqft26: (99626, 0.9990172877140909)
numberofstories: (76935, 0.771199470338687)
taxdelinquencyflag: (46879, 0.974740283181581)
taxdelinquencyyear: (46879, 0.974740283181581)
Length: 26
```

Esistono 26 righe con una percentuale di valori assenti oltre il 70%, molte delle quali sono superiori al 95%.

```
In [39]: # Given a dataframe and some column names returns the dataframe within that columns
def remove_columns(df, col_names):
    df.drop(col_names, axis=1, inplace=True)
    return df
```

```
In [40]: col_to_delete
```

```
Out[40]: ['airconditioningtypeid',
'architecturalstyletypeid',
'basementsqft',
'buildingclastypeid',
'decktypeid',
'finishedfloorisquarefeet',
'finishedsquarefeet12',
'finishedsquarefeet15',
'finishedsquarefeet50',
'finishedsquarefeet150',
'garagecarcnt',
'garagetotalsqft',
'hashottuborspa',
'poolsizeusm',
'pooltypeid0',
'pooltypeid2',
'pooltypeid7',
'regionidneighborhood',
'storytypeid',
'threequarterbathnbr',
'typeconstructiontypeid',
'yardbuildingsqft17',
'yardbuildingsqft26',
'taxdelinquencyflag',
'taxdelinquencyyear']
```

```
In [42]: for X in [X_train, X_val, X_test]:
            X = remove_columns(X, col_to_delete)
```

```
Out[42]: dimensionality()
```

```
Out[43]: X_train (99724, 31)
X_val (33578, 31)
X_test (33578, 31)
```

Le colonne sono state rimosse correttamente.

## Rimozione di Feature ridondanti

Alcune feature portano informazione ripetuta: due colonne diverse contribuiscono con lo stesso tipo di informazione.

### fireplaceflag & fireplacecnt

fireplaceflag e fireplacecnt: la prima spiega se esiste almeno un impianto, la seconda quanti impianti sono presenti. La seconda feature porta una informazione almeno uguale a quella della prima.

```
In [43]: X_train.loc[:, ['fireplacecnt', 'fireplaceflag']].head(20)
```

```
Out[43]:
   fireplacecnt  fireplaceflag
0              0              0.0
1              0              0.0
2              0              0.0
3              0              0.0
4              0              0.0
5              1              1.0
6              0              0.0
7              0              0.0
8              0              0.0
9              0              0.0
10             0              0.0
11             1              1.0
12             0              0.0
13             0              0.0
14             0              0.0
15             0              0.0
16             0              0.0
17             1              1.0
18             0              0.0
19             0              0.0
```

Non c'è coerenza tra le due feature.

```
In [44]: sum((
            (get_col(X_train, 'fireplacecnt') != 1) &
            X_train.loc[:, 'fireplaceflag'] == 0
            ).values.ravel())
```

```
Out[44]: 10773
```

In 10000 osservazioni in cui il flag dice che non ci sono impianti, se ne conta almeno uno.

```
In [45]: sum((
            (get_col(X_train, 'fireplacecnt') != 1) &
            X_train.loc[:, 'fireplaceflag'] == 1
            ).values.ravel())
```

```
Out[45]: 229
```

In 200 osservazioni in cui il flag segnala la presenza di un impianto se ne contano zero.

E in egual misura in 200 case dove non si contano impianti il flag ne segnala la presenza.

Scelgo di mantenere l'informazione portata da fireplaceflag, perché più ricca.

```
In [46]: for X in [X_train, X_val, X_test]:
            X = remove_columns(X, ['fireplaceflag'])
```

```
Out[46]: dimensionality()
```

```
Out[47]: X_train (99724, 32)
X_val (33578, 32)
X_test (33578, 32)
```

### fullbathcnt & bathroomcnt

Entrambe le feature conteggiano il numero di bagni.

```
In [48]: X_train.loc[:, ['fullbathcnt', 'bathroomcnt']]
```

```
Out[48]:
   fullbathcnt  bathroomcnt
0              2              2.0
1              1              1.0
2              2              2.0
3              3              3.0
4              1              1.0
...          ...          ...
100727         3              3.0
100728         2              2.0
100729         2              2.5
100730         2              2.5
100731         1              1.0
```

99724 rows x 2 columns

bathroomcnt: porta un'informazione decimale, infatti la sua descrizione cita: including fractional bathrooms.

```
In [49]: sum((get_col(X_train, 'bathroomcnt') - get_col(X_train, 'fullbathcnt') > 1).values.ravel())
```

```
Out[49]: 11
```

In solo una decina di istanze il dato non ha lo stessa parte intera.

```
In [50]: sum((get_col(X_train, 'bathroomcnt') - get_col(X_train, 'fullbathcnt') > 1.5).values.ravel())
```

```
Out[50]: 4
```

E in solo una è maggiore di 1.5.

Scelgo di mantenere solo la colonna bathroomcnt poiché più ricca.

```
In [51]: for X in [X_train, X_val, X_test]:
            X = remove_columns(X, 'fullbathcnt')
```

```
Out[51]: dimensionality()
```

```
Out[52]: X_train (99724, 31)
X_val (33578, 31)
X_test (33578, 31)
```

### fips & censustractblock

fips e censustractblock: contribuiscono con lo stesso tipo di informazione numerica, semplicemente su scala decimale differente.

```
In [53]: get_col(X_train, ['fips', 'censustractblock'])
```

```
Out[53]:
   fips  censustractblock
0    60590    6.059022e+13
1    60370    6.037294e+13
2    60370    6.037294e+13
3    60590    6.059087e+13
4    60370    6.037302e+13
...    ...    ...
100727    60370    6.037571e+13
100728    60590    6.059089e+13
100729    61110    6.111007e+13
100730    60590    6.059022e+13
100731    60370    6.037530e+13
```

99724 rows x 2 columns

Controllo per quante istanze vale questa relazione.

```
In [54]: equal = []
for
```







```
bathroomcnt
[ 2.5  3.  2.  1.  4.  1.5  5.  3.5  0.  4.5  6.  6.5  5.5  8.
 10.  7.5  7.  15.  9.  11.  18.  13.  12.  8.5] (24)

bedroomcnt
[ 3.  4.  5.  0.  2.  1.  6.  8.  7.  9.  10.  12.  11.  16.  14.  13.] (16)

buildingqualitytypeid
[ nan  8.  7.  4.  6.  10.  9.  1.  5.  11.  12.  3.  2.] (13)

calculatedbathnbr
[ 2.5  3.  2.  1.  4.  1.5  5.  3.5  nan  4.5  6.  6.5  5.5  8.
 10.  7.5  7.  15.  9.  11.  18.  13.  12.  8.5] (24)

fips
[6059, 6037, 6111.] (3)

fireplacecnt
[0. 1. 2. 3. 4. 5.] (6)

heatingorsystemtypeid
[ nan  2.  7.  24.  6.  13.  20.  18.  1.  11.  10.  12.] (12)

poolcnt
[0. 1.] (2)

propertycountylandusecode
['128' '0100' 'rare' '34' '010C' '0101'] (6)

propertylandusetypeid
[261. 266. 269. 247. 265. 246. 275. 267. 260. 248. 263. 31. 264.] (13)

propertyzoningdesc
['rare' 'LAR1' 'LAR3' 'LARD1.5' 'LBRIN' 'LAR8'] (6)

regionidcounty
[1286. 3101. 2061.] (3)

roomcnt
[ 7.  0.  8.  6.  5.  4.  9.  10.  3.  2.  11.  15.  12.  13.  14.  1.] (16)

unitcnt
[ nan  1.  3.  2.  4.  143.  9.  70.  6.  237.  11.] (11)

assessmentyear
[2016. 2015.] (2)
```

Analisi dei tipi di valore e divisione in **categoriali** e **ordinali**: divido questi due tipi poiché andranno processati in maniera differente. Per i primi sarà applicato il **One-Hot Encoding**, non verrà applicato per i secondi: comporterebbe una perdita di informazione.

```
In [95]: # Valori discreti: categorici e ordinali
categorical = ['assessmentyear',
               'fips',
               'heatingorsystemtypeid',
               'poolcnt',
               'propertycountylandusecode',
               'propertylandusetypeid',
               'propertyzoningdesc',
               'regionidcounty'
               ]
ordinal = ['bathroomcnt',
           'bedroomcnt',
           'buildingqualitytypeid',
           'calculatedbathnbr',
           'fireplacecnt',
           'roomcnt',
           'unitcnt'
           ]

Fatta eccezione per 'parcelid', il resto dei dati sono numerici.
```

```
In [96]: # Valori continui: numerici
numeric = list(set(X_train.columns) - set(categorical + ordinal) - {'parcelid'})
```

```
In [97]: X_train[numeric].info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 99709 entries, 0 to 100731
Data columns (total 19 columns):
 #   Column                                     Non-Null Count  Dtype
---  --
 0   living_area_prop                         88883 non-null   float32
 1   structuretaxvaluedollarcnt              99436 non-null   float32
 2   taxamount                               99699 non-null   float32
 3   finishedsquarefeet12                    94874 non-null   float32
 4   period_mean_price                       99709 non-null   float64
 5   calculatedfinishedsquarefeet            99224 non-null   float32
 6   longitude                               99709 non-null   float32
 7   tax_prop                                99436 non-null   float32
 8   tax_ratio                               99698 non-null   float32
 9   regionidzip                             99661 non-null   float32
10   landtaxvaluedollarcnt                  99708 non-null   float32
11   taxvaluedollarcnt                       99708 non-null   float32
12   rawcensustractandblock                  99709 non-null   float32
13   regionidcity                            99709 non-null   float64
14   neighborhood_mean_price                 99700 non-null   float64
15   lotsizesquarefeet                       88847 non-null   float32
16   int_transactiondate                     99709 non-null   float64
17   yearbuilt                               99138 non-null   float32
18   latitude                                99709 non-null   float32
dtypes: float32(16), float64(3)
memory usage: 11.1 MB
```

```
In [98]: dimensionality()

X_train (99709, 35)
X_val (33572, 35)
X_test (33567, 35)
```

```
In [99]: len(numeric) + len(categorical) + len(ordinal) + 1
```

Verifica che i sottoinsiemi numerici, categorici e ordinali costituiscano la totalità delle colonne.

```
In [100]: def dim_check():
          return X_train.shape[1] == len(numeric) + len(categorical) + len(ordinal) + 1 # 1: parcelid
```

```
In [101]: dim_check()

True
```

```
Out[101]: True
```

```
In [102]: print(f'numeric: \n{numeric} \n(len(numeric)) \n')
          print(f'categorical: \n{categorical} \n(len(categorical)) \n')
          print(f'ordinal: \n{ordinal} \n(len(ordinal)) \n')
```

```
numeric:
['living_area_prop', 'structuretaxvaluedollarcnt', 'taxamount', 'finishedsquarefeet12', 'period_mean_price', 'c
alculatedfinishedsquarefeet', 'longitude', 'tax_prop', 'tax_ratio', 'regionidzip', 'landtaxvaluedollarcnt', 'la
ndtaxvaluedollarcnt', 'rawcensustractandblock', 'regionidcity', 'neighborhood_mean_price', 'lotsizesquarefeet', 'in
t_transactiondate', 'yearbuilt', 'latitude'] (19)
```

```
categorical:
['assessmentyear', 'fips', 'heatingorsystemtypeid', 'poolcnt', 'propertycountylandusecode', 'propertylandusety
eid', 'propertyzoningdesc', 'regionidcounty'] (8)
```

```
Ordinal:
['bathroomcnt', 'bedroomcnt', 'buildingqualitytypeid', 'calculatedbathnbr', 'fireplacecnt', 'roomcnt', 'unitcnt
'] (7)
```

## Missing-flag: numerici e ordinali

Analisi della percentuale di missing value in questi tipi di dati per controllare se è sensato inserire le missing-flags.

```
In [103]: # Returns rows with nan-percentage over the cut-off
def over_nan_percentage(columns, cutoff, verbose=False):
    over = []
    for cn in columns:
        col = get_col(X_train, cn)
        perc = get_col_nan_info(col)
        if verbose:
            print(f'({cn}): {perc}')
        if perc > cutoff:
            over.append(cn)
    return over
```

```
In [104]: put_nan_flag = over_nan_percentage(numeric+ordinal, 0.2, verbose=True)
put_nan_flag

living_area_prop: 0.1119847151049554
structuretaxvaluedollarcnt: 0.002737967485382463
taxamount: 0.000100391849818089
finishedsquarefeet12: 0.048491091275612
period_mean_price: 0.0
calculatedfinishedsquarefeet: 0.0048644154690148332
longitude: 0.0
tax_prop: 0.002737967485382463
tax_ratio: 0.00011032013405954979
regionidzip: 0.00048140087655076274
landtaxvaluedollarcnt: 1.002918492814089e-05
taxvaluedollarcnt: 1.002918492814089e-05
rawcensustractandblock: 0.0
regionidcity: 0.01941602020880765
neighborhood_mean_price: 0.028264532686e-05
lotsizesquarefeet: 0.10893700668946633
int_transactiondate: 0.0
yearbuilt: 0.0057266645339684484
latitude: 0.0
bathroomcnt: 0.0
bedroomcnt: 0.0
buildingqualitytypeid: 0.3625656643372213
calculatedbathnbr: 0.010149535147278581
fireplacecnt: 0.0
roomcnt: 0.0
unitcnt: 0.35170345706004474
dtypes: float32(16), float64(3)
memory usage: 11.1 MB
```

I missing value hanno una bassissima percentuale. Solo **buildingqualitytypeid** e **unitcnt** hanno una percentuale superiore al 2%. Solo per queste colonne è aggiunto un missing-flag.

```
In [105]: # Given a dataframe and a column_name adds the missing flag
def add_missing_flag(df, col_name):
    df[col_name+'_na_flag'] = df.loc[:,col_name].isna().astype(int)
    return df
```

```
In [106]: for df in [X_train, X_val, X_test]:
          for cname in put_nan_flag:
              df = add_missing_flag(df, cname)
```

```
In [107]: for cname in put_nan_flag:
          print(X_train.loc[:, [cname, cname+'_na_flag']])

           buildingqualitytypeid  buildingqualitytypeid_na_flag
0                               NaN                             0
1                               8.0                             0
2                               7.0                             1
3                               NaN                             1
4                               4.0                             0
...
100727                          7.0                             1
100728                          NaN                             1
100729                          NaN                             1
100730                          NaN                             1
100731                          7.0                             0

[99709 rows x 2 columns]
```

```
           unitcnt  unitcnt_na_flag
0                NaN                1
1                1.0                0
2                1.0                0
3                NaN                1
4                1.0                1
...
100727            1.0                0
100728            NaN                1
100729            NaN                1
100730            NaN                1
100731            1.0                0

[99709 rows x 2 columns]
```

```
In [108]: dimensionality()

X_train (99709, 37)
X_val (33572, 37)
X_test (33567, 37)
```

L'operazione è avvenuta correttamente.

## Riempimento dei Nan - numerici e ordinali

Vista la bassa percentuale di Nan individuata, questi sono riempiti con la mediana della colonna.

```
In [109]: # Given a dataframe and its column names fill its Nans with the median value of the column for that region
def fill_nan_with_median_same_country(df, col_names, country_ids):
    for country_id in country_ids:
        df_sub = df[df.loc[:, 'regionidcounty'] == country_id]
        df_sub = fill_nan_with_median(df_sub, col_names)
    return df

# Given a dataframe and its column names fill its Nans with the median value of the column
def fill_nan_with_median(df, col_names):
    for col_name in col_names:
        df[col_name] = df[col_name].fillna(get_col(df, col_name).median())
    return df
```

```
In [110]: for X in [X_train, X_val, X_test]:
          X = fill_nan_with_median(X, numeric+ordinal, [1286., 2061., 3101.] )

# X = fill_nan_with_median(X, numeric+ordinal, [1286., 2061., 3101.] )
```

```
In [111]: X_train[numeric + ordinal].info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 99709 entries, 0 to 100731
Data columns (total 26 columns):
 #   Column                                     Non-Null Count  Dtype
---  --
 0   living_area_prop                         99709 non-null   float32
 1   structuretaxvaluedollarcnt              99709 non-null   float32
 2   taxamount                               99709 non-null   float32
 3   finishedsquarefeet12                    99709 non-null   float64
 4   period_mean_price                       99709 non-null   float32
 5   calculatedfinishedsquarefeet            99709 non-null   float32
 6   longitude                               99709 non-null   float32
 7   tax_prop                                99709 non-null   float32
 8   tax_ratio                               99709 non-null   float32
 9   regionidzip                             99709 non-null   float32
10   landtaxvaluedollarcnt                  99709 non-null   float32
11   taxvaluedollarcnt                       99709 non-null   float32
12   rawcensustractandblock                  99709 non-null   float32
13   regionidcity                            99709 non-null   float32
14   neighborhood_mean_price                 99709 non-null   float64
15   lotsizesquarefeet                       99709 non-null   float32
16   int_transactiondate                     99709 non-null   float64
17   yearbuilt                               99709 non-null   float32
18   latitude                                99709 non-null   float32
19   bathroomcnt                             99709 non-null   float32
20   bedroomcnt                             99709 non-null   float32
21   buildingqualitytypeid                   99709 non-null   float32
22   calculatedbathnbr                       99709 non-null   float32
23   fireplacecnt                           99709 non-null   float32
24   roomcnt                                99709 non-null   float32
25   unitcnt                                99709 non-null   float32
dtypes: float32(23), float64(3)
memory usage: 13.8 MB
```

## One-Hot encoding delle variabili categoriali

Il One-Hot encoding è effettuato alla luce dei valori del Train e rassicurato in maniera opaca su Validation e Test.

Nell'eventualità che questi due dataset presentassero un valore inedito, il suo encoding sarebbe una riga di zeri per le colonne considerate.

```
In [112]: categorical

['assessmentyear',
 'fips',
 'heatingorsystemtypeid',
 'poolcnt',
 'propertycountylandusecode',
 'propertylandusetypeid',
 'propertyzoningdesc',
 'regionidcounty']
```

```
In [113]: # Given a train-dataframe, its column-names and a list of dataframes,
# trains a one-hot-encoder to the train
# makes a one-hot-encoding for each dataframe
def one_hot_encoding(df_fit, col_names, dfs):
    oh = OneHotEncoder(sparse=False, handle_unknown='ignore')
    oh.fit(df_fit[col_names])
    for df in dfs:
        encoded = oh.transform(df[col_names])
        for i, col in enumerate(oh.get_feature_names(col_names)):
            df[col] = encoded[i,:].astype(int)
            df.drop(col_names, axis=1, inplace=True)
```

```
In [114]: one_hot_encoding(X_train, categorical, [X_train, X_val, X_test])
```

```
In [115]: dimensionality()

X_train (99709, 76)
X_val (33572, 76)
X_test (33567, 76)
```

Rimuovo colonne che codificano i Nan per One-Hot-Encoding: mantengo righe di soli zeri.

```
In [116]: nan_column = list(filter(re.compile("^_na$"), X_train.columns))
print(nan_column)

['heatingorsystemtypeid_nan']
```

```
In [117]: X_train.columns

Index(['parcelid', 'bathroomcnt', 'bedroomcnt', 'buildingqualitytypeid',
       'calculatedbathnbr', 'calculatedfinishedsquarefeet',
       'finishedsquarefeet12', 'fireplacecnt', 'latitude', 'longitude',
       'lotsizesquarefeet', 'rawcensustractandblock', 'regionidcity',
       'regionidzip', 'roomcnt', 'unitcnt', 'yearbuilt',
       'structuretaxvaluedollarcnt', 'taxvaluedollarcnt',
       'landtaxvaluedollarcnt', 'taxamount', 'int_transactiondate',
       'period_mean_price', 'neighborhood_mean_price', 'living_area_prop',
       'tax_ratio', 'tax_prop', 'buildingqualitytypeid_na_flag',
       'living_area_prop', 'assessmentyear_2015.0', 'assessmentyear_2016.0',
       'unitcnt_na_flag', 'assessmentyear_2015.0', 'assessmentyear_2016.0',
       'fips_6037.0', 'fips_6059.0', 'fips_6111.0',
       'heatingorsystemtypeid_1.0', 'heatingorsystemtypeid_2.0',
       'heatingorsystemtypeid_6.0', 'heatingorsystemtypeid_7.0',
       'heatingorsystemtypeid_10.0', 'heatingorsystemtypeid_11.0',
       'heatingorsystemtypeid_12.0', 'heatingorsystemtypeid_13.0',
       'heatingorsystemtypeid_18.0', 'heatingorsystemtypeid_20.0',
       'heatingorsystemtypeid_24.0', 'heatingorsystemtypeid_nan',
       'poolcnt_0.0', 'poolcnt_1.0', 'propertycountylandusecode_0100',
       'propertycountylandusecode_010C', 'propertycountylandusecode_010C',
       'propertycountylandusecode_122', 'propertycountylandusecode_34',
       'propertycountylandusecode_rare', 'propertylandusetypeid_31.0',
       'propertylandusetypeid_246.0', 'propertylandusetypeid_247.0',
       'propertylandusetypeid_248.0', 'propertylandusetypeid_260.0',
       'propertylandusetypeid_261.0', 'propertylandusetypeid_263.0',
       'propertylandusetypeid_264.0', 'propertylandusetypeid_265.0',
       'propertylandusetypeid_266.0', 'propertylandusetypeid_267.0',
       'propertylandusetypeid_269.0', 'propertylandusetypeid_272.0',
       'propertyzoningdesc_LAR1', 'propertyzoningdesc_LAR3',
       'propertyzoningdesc_LARD1.5', 'propertyzoningdesc_LAR8',
       'propertyzoningdesc_LBRIN', 'propertyzoningdesc_rare',
       'regionidcounty_1286.0', 'regionidcounty_2061.0',
       'regionidcounty_3101.0'],
      dtype='object')
```

E colonne che potrebbero essere binarie, come **poolcnt** o **assessmentyear**.

```
In [118]: for X in [X_train, X_val, X_test]:
          X = remove_column(X, nan_column)
          X = remove_column(X, 'poolcnt_0.0')
          X = remove_column(X, 'assessmentyear_2016.0')
```

```
In [119]: dimensionality()

X_train (99709, 73)
X_val (33572, 73)
X_test (33567, 73)
```

## Scrittura csv

Salvo i dati processati in una specifica cartella.

```
In [120]: dir_name = 'preparazione'

X_train.to_csv(dir_name + '/X_train.csv', index=False)
X_val.to_csv(dir_name + '/X_val.csv', index=False)
X_val.to_csv(dir_name + '/X_val.csv', index=False)
X_test.to_csv(dir_name + '/X_test.csv', index=False)
X_test.to_csv(dir_name + '/X_test.csv', index=False)
```