

# 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 generali, 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
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_v2.csv"
fp.train2017 = dir_name + "/train_2017.csv"
```

```
In [3]: # Importo i dati
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 (366321, 58)  
Train 2016 (90275, 3)  
Properties 2017 (982521, 58)  
Train 2017 (77613, 3)

Non dispongo del log\_error di ogni casa, ma solo di quelle che sono state vendute: seleziono solo l'insieme di case di cui ho a disposizione il log\_error.

Unione in un unico dataset: matengo 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 (90275, 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  int32
 1   airconditioningtypeid                53788 non-null  float64
 2   architecturaltypeid                  468 non-null    float64
 3   basementqft                          93 non-null     float64
 4   bathroomcnt                          167854 non-null float64
 5   bedroomcnt                           167854 non-null float64
 6   buildingqualitytypeid                107173 non-null float64
 7   buildingqualitytypeid                107173 non-null float64
 8   calculatedbathnbr                    166056 non-null float64
 9   decktypeid                           1272 non-null   float64
10  finishedfloorisquarefeet             12893 non-null   float64
11  calculatedfinishedsquarefeet         166992 non-null float64
12  finishedsquarefeet12                  159519 non-null float64
13  finishedsquarefeet13                  125 non-null    float64
14  finishedsquarefeet15                   6591 non-null   float64
15  finishedsquarefeet50                   807 non-null    float64
16  finishedsquarefeet6                   807 non-null    float64
17  fips                                   167854 non-null float64
18  fireplacescnt                         17896 non-null   float64
19  fullbathcnt                           166056 non-null float64
20  garagecarcnt                          55457 non-null   float64
21  garagetotalqft                       55457 non-null   float64
22  hashottuborspa                       3904 non-null    object
23  heatingorsystemtypeid                 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  poolsizeum                             1838 non-null    float64
29  pooltypeid10                          1626 non-null    float64
30  pooltypeid2                            2278 non-null    float64
31  pooltypeid7                            31776 non-null   object
32  propertycountylandusecode             167853 non-null object
33  propertylandusecode                   167854 non-null float64
34  propertyzoningdesc                     108789 non-null object
35  rawcensustrackandblock                 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                 103056 non-null float64
44  unitcnt                               103056 non-null float64
45  yearbuildingqft17                     5039 non-null    float64
46  yearbuildingqft26                     165 non-null     float64
47  yearbuilt                             166828 non-null float64
48  numberofstories                       38169 non-null   object
49  fireplaceflag                         1394 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 float64
55  taxdelinquencyflag                     4683 non-null    float64
56  taxdelinquencyyear                     167002 non-null float64
57  censustrackandblock                   167888 non-null object
58  logerror                              167888 non-null object
59  transactiondate                       167888 non-null object
dtypes: float32(3), int32(1), object(6)
memory usage: 76.94 MB
```

## Casting dei tipi

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

```
In [10]: # Given a dataframe cast all numeric type from 64 bit to 32 bit
def cast_df_to_32(df):
    if dtype == np.float64:
        df[c] = df[c].astype(np.float32)
    if 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  int32
 1   airconditioningtypeid                53788 non-null  float64
 2   architecturaltypeid                  468 non-null    float32
 3   basementqft                          93 non-null     float32
 4   bathroomcnt                          167854 non-null float32
 5   bedroomcnt                           167854 non-null float32
 6   buildinglasttypeid                   31 non-null     float32
 7   buildingqualitytypeid                107173 non-null float32
 8   calculatedbathnbr                    166056 non-null float32
 9   decktypeid                           1272 non-null   float32
10  finishedfloorisquarefeet             12893 non-null   float32
11  calculatedfinishedsquarefeet         166992 non-null float32
12  finishedsquarefeet12                  159519 non-null float32
13  finishedsquarefeet13                  125 non-null    float32
14  finishedsquarefeet15                   6591 non-null   float32
15  finishedsquarefeet50                   807 non-null    float32
16  finishedsquarefeet6                   807 non-null    float32
17  fips                                   167854 non-null float32
18  fireplacescnt                         17896 non-null   float32
19  fullbathcnt                           166056 non-null float32
20  garagecarcnt                          55457 non-null   float32
21  garagetotalqft                       55457 non-null   float32
22  hashottuborspa                       3904 non-null    object
23  heatingorsystemtypeid                 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  poolsizeum                             1838 non-null    float32
29  pooltypeid10                          1626 non-null    float32
30  pooltypeid2                            2278 non-null    float32
31  pooltypeid7                            31776 non-null   object
32  propertycountylandusecode             167853 non-null object
33  propertylandusecode                   167854 non-null float32
34  propertyzoningdesc                     108789 non-null object
35  rawcensustrackandblock                 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                 103056 non-null float32
44  unitcnt                               103056 non-null float32
45  yearbuildingqft17                     5039 non-null    float32
46  yearbuildingqft26                     165 non-null     float32
47  yearbuilt                             166828 non-null float32
48  numberofstories                       38169 non-null   object
49  fireplaceflag                         1394 non-null    object
50  structuretaxvaluedollarcnt            167359 non-null float32
51  taxvaluedollarcnt                     167852 non-null float32
52  assessmentyear                         167854 non-null float32
53  landtaxvaluedollarcnt                 167851 non-null float32
54  taxamount                             167843 non-null float32
55  taxdelinquencyflag                     4683 non-null    float32
56  taxdelinquencyyear                     167002 non-null float32
57  censustrackandblock                   167888 non-null object
58  logerror                              167888 non-null object
59  transactiondate                       167888 non-null object
dtypes: float32(3), 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
11991059    3
11842707    3
14010551    3
12478591    3
14672826    3
17164212    3
17237150    3
12612211    3
11186156    2
11991474    2
12273962    2
11061551    2
14659784    2
12467034    2
11266554    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().index[dfAll.loc[:, 'parcelid'].value_counts() == 3].to_dict().keys)
houses
Out[15]: [10857130, 11991059, 11842707, 14010551, 12478591, 14672826, 17164212, 17237150, 12612211, 11186156, 11991474, 12273962, 11061551, 14659784, 12467034, 11266554, 14008322, 12752161, 11887100, 12239653]
```

Ispeziono logerror e transactiondate di queste case.

```
In [16]: dfAll.loc[:, 'parcelid'].isin(houses)
Out[16]:
parcelid  logerror  transactiondate
135236    10857130    0.053244    2017-06-30
135237    10857130    0.053244    2017-06-30
135238    10857130    0.290908    2017-08-25
55794     11842707    -0.028040    2016-07-14
55795     11842707    0.057900    2016-08-22
55796     11842707    0.207800    2016-09-29
134115    11991059    2.619876    2017-06-06
134116    11991059    2.670239    2017-06-09
134117    11991059    2.508444    2017-06-13
48461     12478591    0.424000    2017-06-23
97365     12478591    0.012482    2017-02-01
97366     12478591    0.039378    2017-09-18
136926    12612211    -0.007561    2017-06-15
136927    12612211    0.074989    2017-08-31
136928    12612211    0.089218    2017-09-19
20217     14010551    1.021000    2016-09-29
114234    14010551    -0.005612    2017-04-06
114235    14010551    0.082468    2017-08-11
33582     14672826    -0.013100    2016-05-10
33583     14672826    0.007000    2017-08-11
155870    14672826    -0.000484    2017-08-11
15385     17164212    -0.080100    2016-09-10
15386     17164212    -0.018200    2016-02-19
165695    17164212    -0.003952    2017-09-11
10416     17237150    0.213500    2016-02-19
10417     17237150    0.288900    2016-07-11
104777    17237150    -0.053883    2017-03-03
```

Si nota che il logerror nella stessa casa varia di molto a seconda della data di vendita; nella preparazione dei dati sarà dunque molto importante prendere in considerazione anche il fattore temporale.

## Split in Train, Validation e Test

Separazione dei dataframe mantenendo in X tutte le colonne fatta eccezione per il logerror, che sarà la unica colonna di y.

```
In [17]: # Given a dataframe and the column-target name, returns two dataframes:
# X with all columns except for the target
# y with the only target column
def split_X_y(df, yname):
    Xnames = list(dfAll.columns)
    Xnames.remove(yname)
    X = df.loc[:, Xnames]
    y = df.loc[:, yname]
    return X, y
```

```
In [18]: df_X, df_y = split_X_y(dfAll, 'logerror')
```

Divisione in Train, Validation e Test con proporzioni 6:2:2

```
In [19]: # Splits the given X and y dataset in three parts:
# - train 0.6
# - validation 0.2
# - test 0.2
def train_validation_test(X, y):
    X_train_80, X_test, y_train_80, y_test = train_test_split(X, y,
                                                                test_size=0.20, random_state=2)
    X_train, X_val, y_train, y_val = train_test_split(X_train_80, y_train_80,
                                                       test_size=0.25, random_state=42)
    return X_train, y_train, X_val, y_val, X_test, y_test
```

```
In [20]: X_train, y_train, X_val, y_val, X_test, y_test = train_validation_test(df_X, df_y)
```

```
In [21]: del(dfAll)
```

Definizione di una funzione che dia informazione sulle dimensionalità dei dataset, che ricorra nel corso delle operazioni per verificare il corretto esito delle trasformazioni impiegate.

```
In [22]: # Prints shape of X_train, X_val and X_test
# If y flag is on, also prints y shapes
def dimensionality(y=False):
    print(f'X_train {X_train.shape}')
    print(f'X_val {X_val.shape}')
    print(f'X_test {X_test.shape}')
    if y:
        print(f'y_train {y_train.shape}')
        print(f'y_val {y_val.shape}')
        print(f'y_test {y_test.shape}')

dimensionality(y=True)
```

```
X_train (100732, 59)
X_val (33578, 59)
X_test (33578, 59)
y_train (100732,)
y_val (33578,)
y_test (33578,)
```

```
In [24]: X_train.loc[:, ['parcelid', 'transactiondate']].head()
```

```
Out[24]:
parcelid  transactiondate
153597    14217523    2017-08-02
146235    11199964    2017-07-11
25650     12627031    2016-04-15
122564     13992985    2017-05-02
84846     12086693    2016-10-13
```

```
In [25]: y_train.head()
```

```
Out[25]:
0    0.057681
1    -0.010315
2    0.020800
3    0.001967
4    0.020200
Name: logerror, dtype: float32
```

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

```
In [26]: # 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 [27]: for df in [X_train, X_val, X_test, y_train, y_val, y_test]:
df = arrange_rows(df)
```

```
In [28]: X_train.loc[:, ['parcelid', 'transactiondate']].head()
```

```
Out[28]:
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 [29]: y_train.head()
```

```
Out[29]:
0    0.057681
1    -0.010315
2    0.020800
3    0.001967
4    0.020200
Name: logerror, dtype: float32
```

## Rappresentazione non corretta dei Nan

Analizzando il dataset di 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 opportuno rappresentarle con zero e con 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.
        df.loc[:, col_name][is_na] = 1.
    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, 'fireplace', '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 X and y, a dataframe remove all rows which target value is under the first or over the last percentile
def X_remove_outlier(X, y):
    out1 = y < np.percentile(y, 99.5)
    out2 = y > np.percentile(y, 0.5)
    out = list(map(lambda o1, o2: o1 and o2, out1, out2))
    X = X[out]
    y = y[out]
    return X, y
```

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

```
X_train (100732, 59)
X_val (33578, 59)
X_test (33578, 59)
```

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

```
In [36]: dimensionality(y=True)
```

```
X_train (99724, 59)
X_val (33578, 59)
X_test (33578, 59)
y_train (99724,)
y_val (33578,)
y_test (33578,)
```

Sono state rimosse circa un migliaio di righe dal Train.

## Rimozione colonne con alta percentuale di Nan

Rimozione delle colonne con un'alta percentuale di valori assenti: queste arricchiscono l'informazione dei 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]

# Given a column returns Nan-count and Nan-percentage
def get_col_nan_info(col):
    count = col.isna().sum()
    tot = len(col)
    perc = count/tot
    return count, perc

# 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, cutoff):
    names = df.columns
    overPercentage = []
    for name in names:
        col = get_col(df, name)
        _, perc = get_col_nan_info(col)
        if perc > cutoff:
            overPercentage.append(name)
    return overPercentage
```

```
In [38]: col_to_delete = get_cols_over_nan_percentage(X_train, 0.6)

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)}')
```

```
airconditioningtypeid: (57655, 0.678224459508242)
architecturaltypeid: (99455, 0.99702550519433)
basementqft: (99676, 0.999186715334233)
buildinglasttypeid: (49709, 0.99829379526451)
decktypeid: (99016, 0.992900405118126)
finishedfloorisquarefeet: (91940, 0.9219445670049377)
finishedsquarefeet12: (98676, 0.989518671534833)
finishedsquarefeet13: (95882, 0.9614736673218082)
finishedsquarefeet15: (91940, 0.9219445670049377)
finishedsquarefeet50: (99264, 0.995387648620593)
garagecarcnt: (66611, 0.6679535518029762)
garagetotalqft: (59904, 0.606979267265232)
hashottuborspa: (97984, 0.976655140186946)
poolsizeum: (98619, 0.9889194175925554)
pooltypeid10: (98711, 0.9889419638201436)
pooltypeid2: (98407, 0.9867935019884849)
pooltypeid7: (80817, 0.8104067225542497)
regionidneighborhood: (59904, 0.606979267265232)
storytypeid: (98676, 0.989518671534833)
threequarterbathnbr: (86493, 0.867328137258934)
typeconstructiontypeid: (99417, 0.9969251033492439)
yearbuildingqft17: (86465, 0.974602918318581)
yearbuildingqft26: (98626, 0.9990172877140909)
numberofstories: (76959, 0.7711994705386877)
taxdelinquencyflag: (86979, 0.974740283181581)
taxdelinquencyyear: (86979, 0.974740283181581)
Length: 26
```

Esistono ben 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', 'architecturaltypeid', 'basementqft', 'buildinglasttypeid', 'decktypeid', 'finishedfloorisquarefeet', 'finishedsquarefeet12', 'finishedsquarefeet13', 'finishedsquarefeet15', 'finishedsquarefeet50', 'finishedsquarefeet6', 'garagecarcnt', 'garagetotalqft', 'hashottuborspa', 'poolsizeum', 'pooltypeid10', 'pooltypeid2', 'pooltypeid7', 'regionidneighborhood', 'storytypeid', 'threequarterbathnbr', 'typeconstructiontypeid', 'yearbuildingqft17', 'yearbuildingqft26', 'numberofstories', 'taxdelinquencyflag', 'taxdelinquencyyear']
```

```
In [41]: for X in [X_train, X_val, X_test]:
X = remove_columns(X, col_to_delete)
```

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

```
X_train (99724, 33)
X_val (33578, 33)
X_test (33578, 33)
```

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 & fireplace

fireplaceflag e fireplace: la prima dice 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[:, ['fireplace', 'fireplaceflag']].head(20)
```

```
Out[43]:
fireplace  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          1              1.0
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, ['fireplace']) > X_train.loc[:, 'fireplaceflag']) == 0 > 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, ['fireplace']) > X_train.loc[:, 'fireplaceflag']) == 1) == 0 .values.ravel())
Out[45]: 229
```

In 200 osservazioni in cui il flag dice che è presente un impianto non se ne contano 6 in egual misura in 200 case dove non si contano impianti il flag ne segnala la presenza.

Scelgo di mantenere l'informazione portata da fireplace: perché più ricca.

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

```
In [47]: dimensionality()
```

```
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            2.0          2.5
1            3.0          3.0
2            2.0          2.0
3            2.0          2.0
4            1.0          1.0
...         ...         ...
100727       2.0          2.0
100728       2.0          2.0
100729       2.0          2.5
100730       1.0          1.0
100731       1.0          1.0
```

99724 rows × 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]: 1
```

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')
```

```
In [52]: dimensionality()
```

```
X_train (99724, 31)
X_val (33578, 31)
X_test (33578, 31)
```

### fips & censustrackblock

fips e censustrackblock contribuiscono con esattamente la stessa informazione numerica, semplicemente su scala esponenziale differente.

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

```
Out[53]:
fips  censustrackandblock
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
...     ...     ...
1
```







```
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 sarò fatto **One-Hot Encoding**, per i secondi invece no poiché 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

Out[99]: 35
```

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()

Out[101]: True

In [102]: print(f'Numeric:\n{numeric}. (len(numeric))\n\n')
          print(f'Categorical:\n{categorical}. (len(categorical))\n\n')
          print(f'Ordinal:\n{ordinal}. (len(ordinal))\n\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

Controllo la 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(colnames, cutoff, verbose=False):
    over = []
    for cn in colnames:
        col = get_col(X_train, cn)
        pct = get_col_nan_info(col)
        if verbose:
            print(f'({cn}): {pct}')
        if pct > 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: 0.000100391849818081
finishedsquarefeet12: 0.048491091275612
period_mean_price: 0.0
calculatedfinishedsquarefeet: 0.004864154690148332
longitude: 0.0
tax_prop: 0.002737967485382463
tax_ratio: 0.00011032013405954979
regionidzip: 0.0004814008765076274
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: 0.0057266643539684484
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: 13.8 MB

Out[104]: 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 [105]: # 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)
          # X = fill_nan_with_median_same_country(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   float32
 4   period_mean_price     99709 non-null   float64
 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 rappresento 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

Out[112]: ['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("^_nan$").match, list(X_train.columns)))
print(nan_column)
['heatingorsystemtypeid_nan']

In [117]: X_train.columns

Out[117]: Index(['parcelid', 'bathroomcnt', 'bedroomcnt', 'buildingqualitytypeid',
          'calculatedbathnbr', 'calculatedfinishedsquarefeet',
          'finishedsquarefeet12', 'fireplacecnt', 'latitude', 'longitude',
          'lotsizesquarefeet', 'rawcensustractandblock', 'regionidcity',
          'regionidzip', 'roomcnt', 'unitcnt', 'yearbuilt',
          'structuretaxvaluedollarcnt', 'taxamount', 'int_transactiondate',
          'landtaxvaluedollarcnt', 'taxvaluedollarcnt',
          'period_mean_price', 'neighborhood_mean_price', 'living_area_prop',
          'tax_ratio', 'tax_prop', 'buildingqualitytypeid_na_flag',
          '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_0101', 'propertycountylandusecode_010C',
          'propertycountylandusecode_122', 'propertycountylandusecode_34',
          'propertycountylandusecode_rare', 'propertylandusetypeid_31.0',
          'propertylandusetypeid_246.0', 'propertylandusetypeid_247.0',
          'propertylandusetypeid_248.0', 'propertylandusetypeid_265.0',
          '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_LBRN', 'propertyzoningdesc_rare',
          'regionidcounty_1286.0', 'regionidcounty_2061.0',
          'regionidcounty_3101.0'],
          dtype='object')
```

E colonne che potrebbero essere binarie, come poolcnt o assesmentyear.

```
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_test.to_csv(dir_name + '/X_test.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 + '/y_test.csv', index=False)
```