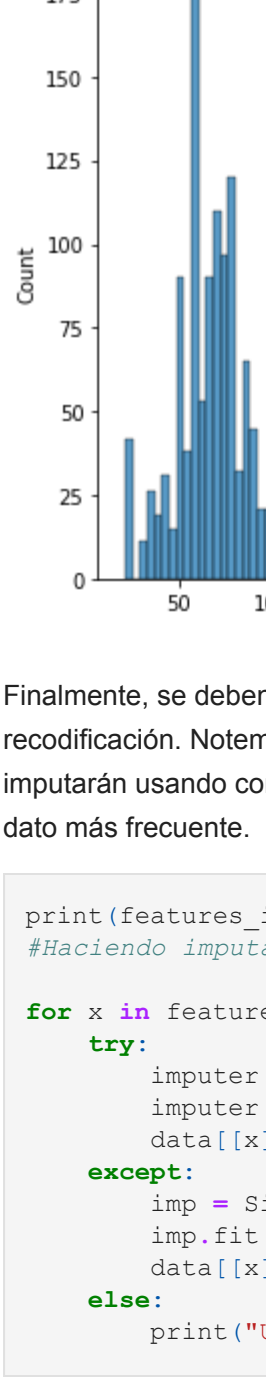


Histogramas después de imputación
 (2919, 81)



Finalmente, se deben hacer imputaciones sobre el resto de las variables que tienen NaNs y que no se pudieron arreglar mediante redefinición. Notemos que algunas de ellas son numéricas y otras son categóricas. De esa cuenta, en el caso de las numéricas, se imputarán usando como dato de imputación la media aritmética, mientras que para las categóricas se usará como dato de imputación el dato más frecuente.

```
In [364]: print(features_imputation)
#Haciendo imputación

for x in features_imputation:
    try:
        imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
        imputer = imputer.fit(data[[x]])
        data[x] = imputer.transform(data[[x]])
    except:
        imp = SimpleImputer(missing_values=np.nan, strategy='most_frequent')
        imp.fit(data[[x]])
        data[x] = imp.transform(data[[x]])
    else:
        print("Un error ha ocurrido mientras se imputaba", x)

['MSZoning', 'Utilities', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'MasVnrArea', 'BsmQual', 'BsmCond', 'BsmtExposure', 'BsmFinType1', 'BsmFinFinS1', 'BsmFinType2', 'BsmFinFinS2', 'BsmUnfSF', 'TotalBsmSF', 'Electrica
l', 'BsmFullBath', 'BsmHalfBath', 'KitchenQual', 'Functional', 'GarageType', 'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual', 'GarageCond', 'SaleType']
Un error ha ocurrido mientras se imputaba MasVnrArea
Un error ha ocurrido mientras se imputaba BsmFinS1
Un error ha ocurrido mientras se imputaba BsmFinS2
Un error ha ocurrido mientras se imputaba BsmUnfSF
Un error ha ocurrido mientras se imputaba TotalBsmSF
Un error ha ocurrido mientras se imputaba BsmFullBath
Un error ha ocurrido mientras se imputaba BsmHalfBath
Un error ha ocurrido mientras se imputaba GarageYrBlt
Un error ha ocurrido mientras se imputaba GarageCars
Un error ha ocurrido mientras se imputaba GarageArea

In [365]: print(data.shape)
print(data.FireplaceQu == data.FireplaceQu.fillna("NS"))
data.Fence = data.Fence.fillna("NS")
data.Alley = data.Alley.fillna("NoAlley")
data.MiscFeature = data.MiscFeature.fillna("NoMiscFeature")
data.PoolQC = data.PoolQC.fillna("NoP")

print(data.shape)
```

(2919, 81)
 (2919, 81)

Ahora se atienden las variables que tienen bajo porcentaje de NaNs, haciendo imputaciones para el resto de las variables (cuyos porcentajes de NA fueron menor a 10%).

```
In [366]: features_imputation = []

for x in porcentajes:
    if porcentajes[x] <10 and porcentajes[x] >0:
        features_imputation.append(x)
        print(x, porcentajes[x])
print(features_imputation)

data.MasVnrType = data.MasVnrType.fillna("NoMasozny")
data.MasVnrArea = data.MasVnrArea.fillna(0)
data.BsmQual = data.BsmQual.fillna("NS") #Cambiar a nominal
data.BsmCond = data.BsmCond.fillna("NS") #Cambiar a nominal
data.BsmExposure = data.BsmExposure.fillna("NM") #Cambiar a nominales
data.BsmFinType1 = data.BsmFinType1.fillna("NB") #Cambiar a nominales
data.BsmFinType2 = data.BsmFinType2.fillna("NB") #Cambiar a nominales
#Con electrical el dato si está corrido, se quiere perder esos registros
data.GarageType = data.GarageType.fillna("NO")

#Imputación
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
imputer = imputer.fit(data[["GarageYrBlt"]])
data.GarageYrBlt = imputer.transform(data[["GarageYrBlt"]])

#...

data.GarageFinish = data.GarageFinish.fillna("NO") #Cambiar a nominal si necesario
data.GarageQual = data.GarageQual.fillna("NO") #Cambiar a nominales
data.GarageCond = data.GarageCond.fillna("NO") #Cambiar a nominal

porcentajes = {}

for column in list(data.columns):
    porcentajes[column] = len(data[column][data[column].isna()==True]) / len(data[column]) *100
print(sorted(porcentajes.items(), key=lambda x:x[1]))
drop = []
print("Porcentaje vale", porcentajes)
for x in porcentajes.keys():
    if porcentajes[x] >0 and x!="SalePrice":
        drop.append(x)
        data.loc["train"].dropna(subset=[x], inplace=True)
print(x)
print(data.shape)

MSZoning 0.137033230584104
Utilities 0.068166152792052
Exterior1st 0.0342583076396026
Exterior2nd 0.0342583076396026
MasVnrType 0.8221993833504623
MasVnrArea 0.78794107570318599
BsmQual 2.774922918807811
BsmCond 2.8091812264474134
BsmExposure 2.8091812264474134
BsmFinType1 2.70640303286057
BsmFinFinS1 0.0342583076396026
BsmFinType2 2.740664611182083
BsmFinFinS2 0.0342583076396026
BsmUnfSF 0.0342583076396026
TotalBsmSF 0.0342583076396026
Electrical 0.0342583076396026
BsmFullBath 0.068166152792052
BsmHalfBath 0.068166152792052
KitchenQual 0.0342583076396026
Functional 0.068166152792052
GarageType 5.378554299417608
GarageYrBlt 5.4470709146968135
GarageFinish 5.4470709146968135
GarageQual 0.0342583076396026
GarageArea 0.0342583076396026
GarageCars 5.4470709146968135
GarageCond 5.4470709146968135
SaleType 0.0342583076396026

['MSZoning', 'Utilities', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'MasVnrArea', 'BsmQual', 'BsmCond', 'BsmtExposure', 'BsmFinType1', 'BsmFinFinS1', 'BsmFinType2', 'BsmFinFinS2', 'BsmUnfSF', 'TotalBsmSF', 'Electrica
l', 'BsmFullBath', 'BsmHalfBath', 'KitchenQual', 'Functional', 'GarageType', 'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual', 'GarageCond', 'SaleType']
[('Id', 0.0), ('MSSubClass', 0.0), ('MSZoning', 0.0), ('LotFrontage', 0.0), ('LotArea', 0.0), ('Street', 0.0), ('Alley', 0.0), ('LotShape', 0.0), ('LandContour', 0.0), ('Utilities', 0.0), ('LotConfig', 0.0), ('LandSlope', 0.0), ('Neighborhood', 0.0), ('Condition', 0.0), ('Condition2', 0.0), ('BldgType', 0.0), ('HouseStyle', 0.0), ('OverallQual', 0.0), ('OverallCond', 0.0), ('YearBuilt', 0.0), ('YearRemodAdd', 0.0), ('RoofStyle', 0.0), ('RoofMatl', 0.0), ('Exterior1st', 0.0), ('Exterior2nd', 0.0), ('ExteriorSf', 0.0), ('CentralAir', 0.0), ('MasVnrType', 0.0), ('MasVnrArea', 0.0), ('BsmQual', 0.0), ('BsmCond', 0.0), ('BsmExposure', 0.0), ('BsmFinType1', 0.0), ('BsmFinFinS1', 0.0), ('BsmFinType2', 0.0), ('BsmFinFinS2', 0.0), ('BsmUnfSF', 0.0), ('TotalBsmSF', 0.0), ('Heating', 0.0), ('HeatingQC', 0.0), ('CentralAir', 0.0), ('Electrical', 0.0), ('FullBath', 0.0), ('HalfBath', 0.0), ('BedroomAbvGr', 0.0), ('KitchenAbvGr', 0.0), ('KitchenQual', 0.0), ('TotRmsAbvGrd', 0.0), ('Functional', 0.0), ('Fireplaces', 0.0), ('FireplaceQu', 0.0), ('GarageType', 0.0), ('GarageYrBlt', 0.0), ('GarageFinish', 0.0), ('GarageCars', 0.0), ('GarageArea', 0.0), ('GarageQual', 0.0), ('GarageCond', 0.0), ('PavedDrive', 0.0), ('WoodDeckSF', 0.0), ('OpenPorchSF', 0.0), ('EnclosedPorch', 0.0), ('3SsnPorch', 0.0), ('ScreenPorch', 0.0), ('PoolArea', 0.0), ('PoolQC', 0.0), ('Fence', 0.0), ('MiscFeature', 0.0), ('3SsnPorch', 0.0), ('ScreenPorch', 0.0), ('PoolArea', 0.0), ('PoolQC', 0.0), ('Fence', 0.0), ('MiscFeature', 0.0), ('MoSold', 0.0), ('YrSold', 0.0), ('SaleType', 0.0), ('SaleCondition', 0.0), ('SalePrice', 49.3820841802)]

Porcentaje vale {'Id': 0.0, 'MSSubClass': 0.0, 'MSZoning': 0.0, 'LotFrontage': 0.0, 'LotArea': 0.0, 'Street': 0.0, 'Alley': 0.0, 'LotShape': 0.0, 'LandContour': 0.0, 'Utilities': 0.0, 'LotConfig': 0.0, 'LandSlope': 0.0, 'Neighborhood': 0.0, 'Condition': 0.0, 'Condition2': 0.0, 'BldgType': 0.0, 'HouseStyle': 0.0, 'OverallQual': 0.0, 'OverallCond': 0.0, 'YearBuilt': 0.0, 'YearRemodAdd': 0.0, 'RoofStyle': 0.0, 'RoofMatl': 0.0, 'Exterior1st': 0.0, 'Exterior2nd': 0.0, 'ExteriorSf': 0.0, 'CentralAir': 0.0, 'MasVnrType': 0.0, 'MasVnrArea': 0.0, 'BsmQual': 0.0, 'BsmCond': 0.0, 'BsmExposure': 0.0, 'BsmFinType1': 0.0, 'BsmFinFinS1': 0.0, 'BsmFinType2': 0.0, 'BsmFinFinS2': 0.0, 'BsmUnfSF': 0.0, 'TotalBsmSF': 0.0, 'Heating': 0.0, 'HeatingQC': 0.0, 'CentralAir': 0.0, 'Electrical': 0.0, 'FullBath': 0.0, 'HalfBath': 0.0, 'BedroomAbvGr': 0.0, 'KitchenAbvGr': 0.0, 'KitchenQual': 0.0, 'TotRmsAbvGrd': 0.0, 'Functional': 0.0, 'Fireplaces': 0.0, 'FireplaceQu': 0.0, 'GarageType': 0.0, 'GarageYrBlt': 0.0, 'GarageFinish': 0.0, 'GarageCars': 0.0, 'GarageArea': 0.0, 'GarageQual': 0.0, 'GarageCond': 0.0, 'PavedDrive': 0.0, 'WoodDeckSF': 0.0, 'OpenPorchSF': 0.0, 'EnclosedPorch': 0.0, '3SsnPorch': 0.0, 'ScreenPorch': 0.0, 'PoolArea': 0.0, 'PoolQC': 0.0, 'Fence': 0.0, 'MiscFeature': 0.0, 'MoSold': 0.0, 'YrSold': 0.0, 'SaleType': 0.0, 'SaleCondition': 0.0, 'SalePrice': 49.3820841802}
```

```
In [367]: data.info()

<class 'pandas.core.frame.DataFrame'>
MultiIndex: 2919 entries, 'train', 0) to ('test', 1458)
Data columns (total 81 columns):
 #   Column              Non-Null Count  Dtype
---  ---
 0   Id                  2919 non-null    int64
 1   MSSubClass          2919 non-null    object
 2   MSZoning            2919 non-null    object
 3   LotFrontage         2919 non-null    float64
 4   LotArea             2919 non-null    object
 5   Street              2919 non-null    object
 6   Alley              2919 non-null    object
 7   LandContour         2919 non-null    object
 8   Utilities           2919 non-null    object
 9   LotConfig           2919 non-null    object
10   LandSlope           2919 non-null    object
11   Neighborhood         2919 non-null    object
12   Condition            2919 non-null    object
13   Condition2          2919 non-null    object
14   BldgType             2919 non-null    object
15   HouseStyle          2919 non-null    object
16   OverallQual          2919 non-null    int64
17   OverallCond          2919 non-null    int64
18   YearBuilt            2919 non-null    int64
19   YearRemodAdd         2919 non-null    int64
20   RoofStyle           2919 non-null    object
21   RoofMatl            2919 non-null    object
22   Exterior1st          2919 non-null    object
23   Exterior2nd         2919 non-null    object
24   MasVnrType          2919 non-null    object
25   MasVnrArea          2919 non-null    float64
26   BsmQual             2919 non-null    object
27   BsmCond             2919 non-null    object
28   BsmExposure         2919 non-null    object
29   BsmFinType1         2919 non-null    object
30   BsmFinType2         2919 non-null    float64
31   BsmUnfSF            2919 non-null    float64
32   TotalBsmSF          2919 non-null    float64
33   Heating             2919 non-null    object
34   HeatingQC           2919 non-null    object
35   CentralAir          2919 non-null    object
42   Electrical          2919 non-null    object
43   FullBath            2919 non-null    int64
44   HalfBath            2919 non-null    int64
45   LowQualFinSF        2919 non-null    int64
46   GrLivArea           2919 non-null    float64
47   BsmFullBath         2919 non-null    float64
48   BsmHalfBath         2919 non-null    float64
49   FullBath            2919 non-null    int64
50   HalfBath            2919 non-null    int64
51   BedroomAbvGr       2919 non-null    int64
52   KitchenAbvGr       2919 non-null    int64
53   KitchenQual         2919 non-null    object
54   TotRmsAbvGrd       2919 non-null    int64
55   Functional          2919 non-null    object
56   Fireplaces          2919 non-null    int64
57   FireplaceQu         2919 non-null    object
58   GarageType          2919 non-null    object
59   GarageYrBlt         2919 non-null    float64
60   GarageFinish        2919 non-null    object
61   GarageCars          2919 non-null    object
62   GarageArea          2919 non-null    float64
63   GarageQual          2919 non-null    object
64   GarageCond          2919 non-null    object
65   PavedDrive          2919 non-null    object
66   WoodDeckSF         2919 non-null    int64
67   OpenPorchSF        2919 non-null    int64
68   EnclosedPorch       2919 non-null    int64
69   3SsnPorch           2919 non-null    int64
70   ScreenPorch         2919 non-null    int64
71   PoolArea           2919 non-null    object
72   PoolQC              2919 non-null    object
73   Fence              2919 non-null    object
74   MiscFeature         2919 non-null    object
75   MiscVal            2919 non-null    int64
76   MoSold             2919 non-null    object
77   YrSold             2919 non-null    int64
78   SaleType            2919 non-null    object
79   SaleCondition       2919 non-null    object
80   SalePrice           1460 non-null    float64
dtypes: float64(12), int64(24), object(45)
memory usage: 1.9+ MB
```

2. Explorando los datos

En esta parte se hará una exploración de los datos para ganar un poco de intuición a cerca de como se comportan las variables.

```
In [368]: data.describe()

Out[368]:
```

	Id	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmFinFinS1	BsmFin
count	2919.000000	2919.000000	2919.000000	2919.000000	2919.000000	2919.000000	2919.000000	2919.000000	2919.000000	2919.00
mean	1460.000000	69.305795	10168.114080	6.080977	5.564577	1971.312778	1984.284474	102.201312	441.423235	49.58
std	842.787043	21.312345	7886.996359	1.409947	1.113131	30.291442	20.894344	178.620089	455.532750	169.17
min	1.000000	1.000000	1300.000000	1.000000	1.000000	1872.000000	1950.000000	0.000000	0.000000	0.00
25%	730.500000	60.000000	7478.000000	5.000000	5.000000	1953.500000	1965.000000	0.000000	0.000000	0.00
50%	1460.000000	69.305795	9453.000000	6.000000	5.000000	1973.000000	1993.000000	0.000000	369.000000	0.00
75%	2189.500000	78.000000	10570.000000	7.000000	6.000000	2001.000000	2004.000000	163.500000	733.000000	0.00
max	2919.000000	313.000000	215245.000000	10.000000	9.000000	2010.000000	2010.000000	1600.000000	5644.000000	1528.00

8 rows x 36 columns

```
In [369]: print(data.shape)

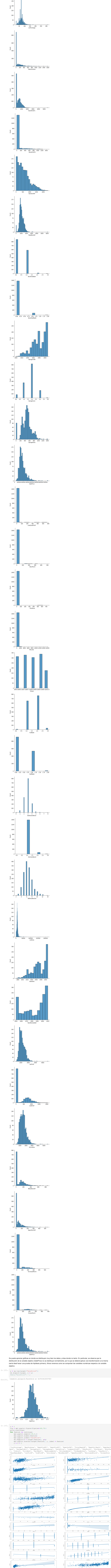
target = data.loc["train"].copy()
target["SalePrice"] = y
#target = target.drop(['SalePrice'], axis=1)
```

Construyendo histogramas para las variables numéricas para ver su comportamiento

```
In [370]: for x in continuas:
sb.displot(target, x = x)
sb.displot(target, x = y)
```


home/hugog/GitHub/yyecto2---Inteligencia-Artificial/venv/lib/python3.9/site-packages/seaborn/axisgrid.py:40: RuntimeWarning: More than 20 figures have been opened. Figures created through the pyplot interface (matplotlib.pyplot.figure) are retained until explicitly closed and may consume too much memory. To control this warning, see the rcParam figure_max_open_warning'. (To control this warning, see the rcParam figure_max_open_warning').

fig = plt.figure(figsize=(8,8))
seaborn.axisgrid.FacetGrid at 0x7f15430f3550>



Se puede apreciar graficas en donde se distribuyen muy bien los datos y otras donde no tanto. En particular, se observa que la distribución de la variable objetivo (SalePrice) no se distribuye normalmente, por lo que se deberá aplicar una transformación a la misma (sería ideal hacer una prueba de hipótesis primero). Ahora veremos como se comportan las variables numéricas respecto a la variable objetivo.

In [37]:
y = np.log(target("SalePrice"))
sb.displot(target, x = y)
data["SalePrice"] = y

Out[37]:
<seaborn.axisgrid.FacetGrid at 0x7f1541f03760>

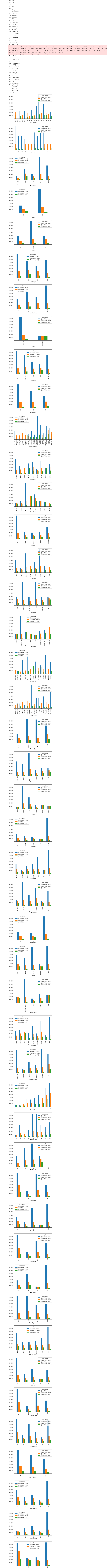


Con estas graficas vemos que casi todas las variables numéricas parecen ser buenas predictoras, salvo:

- BsmFinSF2
- YrSold
- PoolArea

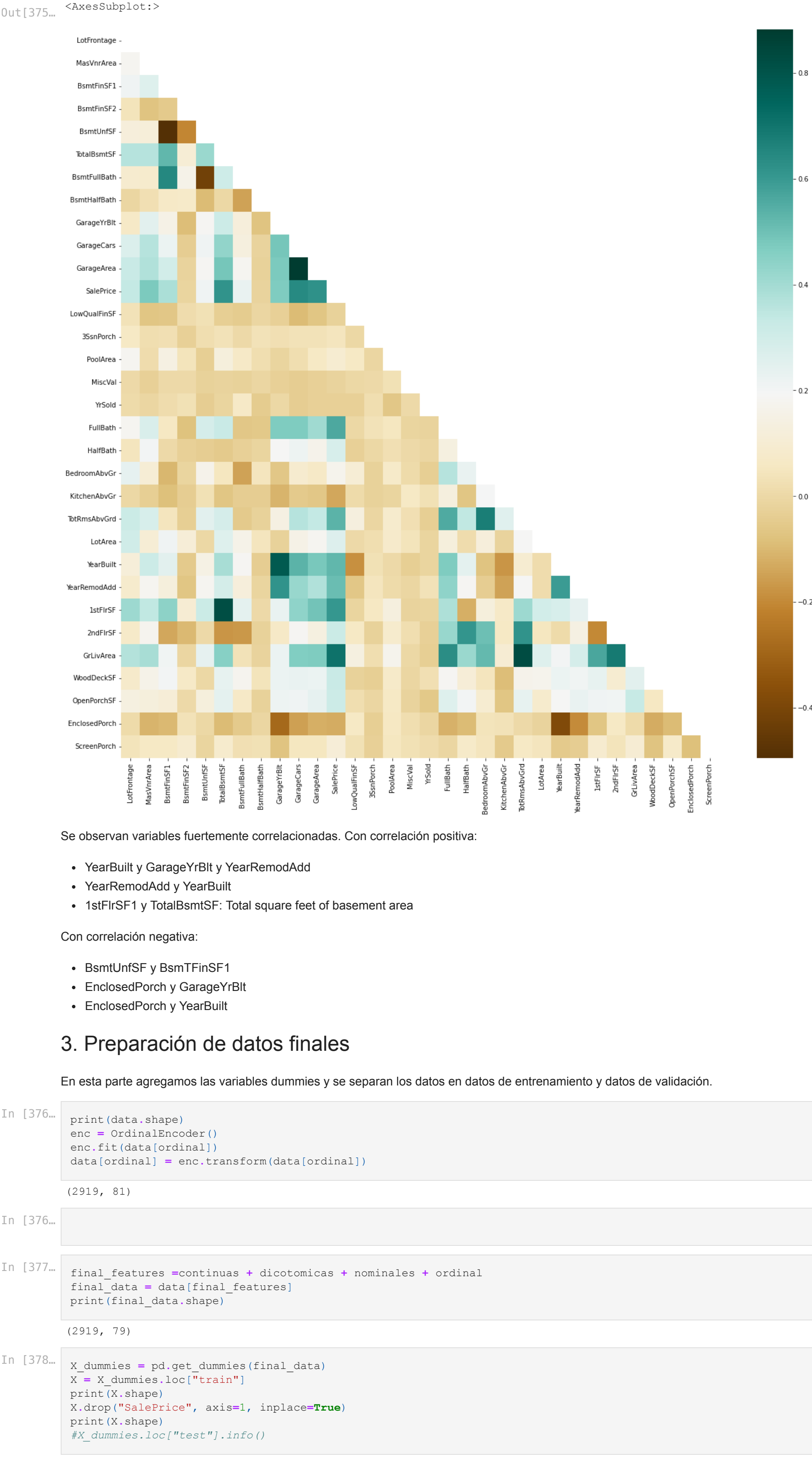
Ahora vemos como lucen las variables categóricas.

In [37]:
for feature in nominales:
 print(feature)
 data[sr(feature), 'SalePrice'].groupby(feature).agg(['mean', 'max', 'min']).plot.bar()
for feature in ordinales:
 print(feature)
 data[sr(feature), 'SalePrice'].groupby(feature).agg(['mean', 'max', 'min']).plot.bar()



De acá vemos, que además de que la variable MoSold ya tenía poco sentido, porque no tener mayor relación con el modelo ya que no se ve diferencias significativas entre las distintas categorías. Esta variable es candidata a ser removida del modelo.

Como es era esperado, la matriz de correlaciones es de 23x23. Como resulta difícil visualizar la matriz de correlaciones se construye una visualización utilizando la biblioteca seaborn.



Se observan variables fuertemente correlacionadas. Con correlación positiva:

- YearBuilt y GarageYrBlt y YearRemodAdd
- YearRemodAdd y YearBuilt
- 1stFlrSF1 y TotalBsmntSF1: Total square feet of basement area

Con correlación negativa:

- BsmUnfSF y BsmTFinSF1
- EnclosedPorch y GarageYrBlt
- EnclosedPorch y YearBuilt

3. Preparación de datos finales

En esta parte agregamos las variables dummies y se separan los datos en datos de entrenamiento y datos de validación.

```
In [376]: print(data.shape)
enc = OrdinalEncoder()
enc.fit(data[ordinal])
data[ordinal] = enc.transform(data[ordinal])

(2919, 81)
```

```
In [376]:
```

```
In [377]: final_features = ['continuos' + dicotomicas + nominales + ordinal
final_data = data[final_features]
print(final_data.shape)

(2919, 79)
```

```
In [378]: X_dummies = pd.get_dummies(final_data)
X = X_dummies.loc["train"]
print(X.shape)
X.drop("SalePrice", axis=1, inplace=True)
print(X.shape)
X_dummies.loc["test"].info()

(1460, 268)
/home/hugob/githu/Projecto2---Inteligencia-Artificial/venv/lib/python3.9/site-packages/pandas/core/frame.py:49
06: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
return super().drop()

```
In [379]: X_train, X_holdout, y_train, y_holdout = train_test_split(X,y, test_size=0.2, random_state=5)
```

4. Modelos propuestos

Ahora haremos una regresión lineal normal utilizando las 269 variables (incluidas las dummies). Cabe destacar que, siguiendo los lineamientos de Gerón, usaremos como medida del desempeño el RSME. En la literatura se acostumbra decir que con un RSME <=0.5 se considera que el modelo tiene buena capacidad de predicción y entre más pequeño sea el RSME mejor es el modelo.

```
In [380]: from math import exp
import statsmodels.api as sm
reg = LinearRegression()
reg.fit(X_train, y_train)

print("*****Estadísticas training*****")
y_predicted = reg.predict(X_train)
estadisticas(y_predicted, y_train)

print("*****Estadísticas testeo*****")
y_predicted = reg.predict(X_holdout)
estadisticas(y_predicted, y_holdout)

*****Estadísticas training*****
R-Square Value 0.9439182737159211

mean_absolute_error : 0.06549455740507339

mean_squared_error : 0.008817315954081558

root_mean_squared_error : 0.09390056418404298
*****Estadísticas testeo*****
R-Square Value -308504.2356959599

mean_absolute_error : 0.1752005494240462

mean_squared_error : 51925.77000896375

root_mean_squared_error : 227.8722668710779
```

```
In [381]: X2 = sm.add_constant(X_train)
est = sm.OLS(y_train,X2)
est2 = est.fit()
```

Con esto vemos que el modelo se sobrentrena, ya que en la parte de entrenamiento logra buen ajuste (R2=0.944) y buen RSME (0.09), pero cuando hacemos la validación con los datos de validación, vemos que el ajuste se pierde (R2=-308404) y RSME demasiado alto (227.9). Esto nos está diciendo que este modelo es malo prediciendo, y sacar el promedio de los datos sería un mejor modelo que esta regresión.

```
In [381]:
```

Debido a que se está observando un sobreentrenamiento (posiblemente al gran número de variables), se hace necesario recurrir a la regularización. Haremos la prueba primero con la regularización de Ridge. Para ello probaremos con valores de alfa muy pequeños y otros más grandes (hasta 100) y compararemos los resultados del RSME.

```
In [382]: from sklearn.preprocessing import StandardScaler
alpha_range = [0.1e-15, 1e-10, 1e-8, 1e-4, 1e-3, 1e-2, 1, 5, 10, 20, 100]
for i in range(len(alpha_range)):
    print("*****Alpha = %.2e*****" % alpha_range[i])
    ridge = make_pipeline(StandardScaler(with_mean=False), Ridge(alpha=alpha_range[i]))
    ridge.fit(X_train, y_train)
    print("*****Estadísticas training Alpha = %.2e*****" % alpha_range[i])
    y_predicted = ridge.predict(X_train)
    estadisticas(y_predicted, y_train)

    print("*****Estadísticas testeo Alpha = %.2e*****" % alpha_range[i])
    y_predicted = ridge.predict(X_holdout)
    estadisticas(y_predicted, y_holdout)

from IPython.display import Image

*****Alpha = 0.00e+00*****
*****Estadísticas training Alpha = 0.00e+00*****
R-Square Value 0.8567693477066085

mean_absolute_error : 0.1166799253319561

mean_squared_error : 0.022518295401031178

root_mean_squared_error : 0.15060697227804163
*****Estadísticas testeo Alpha = 0.00e+00*****
R-Square Value -3.9906151164084527e+27

mean_absolute_error : 156013386656.9604

mean_squared_error : 6.716766648490138e+26

root_mean_squared_error : 25916725263987.613
*****Estadísticas training Alpha = 1.00e-15*****
R-Square Value 0.9439182710593352

mean_absolute_error : 0.06549574390975025

mean_squared_error : 0.00881700193824165

root_mean_squared_error : 0.09398889210337708
*****Estadísticas testeo Alpha = 1.00e-15*****
R-Square Value -3.73693965446496

mean_absolute_error : 0.14751985054209732

mean_squared_error : 0.7973022691883245

root_mean_squared_error : 0.8929180752948864
*****Alpha = 1.00e-10*****
*****Estadísticas training Alpha = 1.00e-10*****
R-Square Value 0.9439182710593353

mean_absolute_error : 0.06549574390970479

mean_squared_error : 0.00881700193824164

root_mean_squared_error : 0.09398889210337704
*****Estadísticas testeo Alpha = 1.00e-10*****
R-Square Value 0.9019789065973366

mean_absolute_error : 0.08628343062000196

mean_squared_error : 0.01649832859585375

root_mean_squared_error : 0.12844581369006913
*****Alpha = 1.00e-08*****
*****Estadísticas training Alpha = 1.00e-08*****
R-Square Value 0.9439182710593352

mean_absolute_error : 0.06549574390975421

mean_squared_error : 0.00881700193824165

root_mean_squared_error : 0.09398889210337705
*****Estadísticas testeo Alpha = 1.00e-08*****
R-Square Value 0.9019789065973366

mean_absolute_error : 0.08628465551928587

mean_squared_error : 0.016498381175788338

root_mean_squared_error : 0.128444602436739075
*****Alpha = 1.00e-04*****
*****Estadísticas training Alpha = 1.00e-04*****
R-Square Value 0.9439182710593359

mean_absolute_error : 0.0654957440886843

mean_squared_error : 0.00881700193824407

root_mean_squared_error : 0.09398889210338996
*****Estadísticas testeo Alpha = 1.00e-04*****
R-Square Value 0.9019789065973366

mean_absolute_error : 0.08628464178149801

mean_squared_error : 0.016498337789381572

root_mean_squared_error : 0.1284458547767353
*****Alpha = 1.00e-03*****
*****Estadísticas training Alpha = 1.00e-03*****
R-Square Value 0.9439182710577867

mean_absolute_error : 0.06549574890116157

mean_squared_error : 0.00881700193848479

root_mean_squared_error : 0.09398889210467012
*****Estadísticas testeo Alpha = 1.00e-03*****
R-Square Value 0.9019789065973366

mean_absolute_error : 0.08628458778938436

mean_squared_error : 0.016498296811769488

root_mean_squared_error : 0.12844569596436264
*****Alpha = 1.00e-02*****
*****Estadísticas training Alpha = 1.00e-02*****
R-Square Value 0.9439182709049682

mean_absolute_error : 0.06549573805592488

mean_squared_error : 0.008817001962510777

root_mean_squared_error : 0.09398889223260718
*****Estadísticas testeo Alpha = 1.00e-02*****
R-Square Value 0.9019815295234496

mean_absolute_error : 0.08628404784972354

mean_squared_error : 0.01649788712049915

root_mean_squared_error : 0.12844410115104513
*****Alpha = 1.00e+01*****
*****Estadísticas training Alpha = 1.00e+01*****
R-Square Value 0.9439168174297813

mean_absolute_error : 0.06550229320809607

mean_squared_error : 0.008817230473538926

root_mean_squared_error : 0.09390010901771587
*****Estadísticas testeo Alpha = 1.00e+01*****
R-Square Value 0.90245991798019

mean_absolute_error : 0.08622485011711884

mean_squared_error : 0.016454215970539555

root_mean_squared_error : 0.12827398789541689
*****Alpha = 5.00e+00*****
*****Estadísticas training Alpha = 5.00e+00*****
R-Square Value 0.943888616468459

mean_absolute_error : 0.06554826500616642

mean_squared_error : 0.008821664301927604

root_mean_squared_error : 0.09392371533285725
*****Estadísticas testeo Alpha = 5.00e+00*****
R-Square Value 0.9031452032702727

mean_absolute_error : 0.08599213553748818

mean_squared_error : 0.01630202466291263

root_mean_squared_error : 0.12767938151053454
*****Alpha = 1.00e+01*****
*****Estadísticas training Alpha = 1.00e+01*****
R-Square Value 0.9438208691958502

mean_absolute_error : 0.06563123822570238

mean_squared_error : 0.008832315168332057

root_mean_squared_error : 0.0939803977868792
*****Estadísticas testeo Alpha = 1.00e+01*****
R-Square Value 0.9040281460998804

mean_absolute_error : 0.08573153013129467

mean_squared_error : 0.016153412782062164

root_mean_squared_error : 0.1270960769733276
*****Alpha = 2.00e+02*****
*****Estadísticas training Alpha = 2.00e+02*****
R-Square Value 0.943827932617562

mean_absolute_error : 0.06584572202465652

mean_squared_error : 0.008862648080973265

root_mean_squared_error : 0.09414163840178938
*****Estadísticas testeo Alpha = 2.00e+02*****
R-Square Value 0.902666708544862

mean_absolute_error : 0.08536424116157658

mean_squared_error : 0.0159444951632004788

root_mean_squared_error : 0.1262733211411056
*****Alpha = 1.00e+02*****
*****Estadísticas training Alpha = 1.00e+02*****
R-Square Value 0.941438555977063

mean_absolute_error : 0.06750905576562441

mean_squared_error : 0.009207437100499353

root_mean_squared_error : 0.0959553912008041
*****Estadísticas testeo Alpha = 1.00e+02*****
R-Square Value 0.9077834704750749

mean_absolute_error : 0.0845536596295811

mean_squared_error : 0.015521338859402786

root_mean_squared_error : 0.12458466542637896

Acá, vemos que el mejor RSME (0.1245) se obtiene con alfa = 100. Hacemos la predicción para los datos de testeo con este modelo y lo subimos a Kaggle. Vemos que obtiene un RSME de 0.14481
```

salida a csv	0.14481
an hour ago by Hugo García	
Prueba 1	

```
In [383]: test = X_dummies.loc["test"]
test.drop("SalePrice", axis=1, inplace=True)
y_predicted = ridge.predict(test)
y_predicted = np.exp(y_predicted)
Id = pd.read_csv("../Datos/test.csv").Id
salida = pd.DataFrame(y_predicted, Id, Id)
salida.to_csv("../Datos/salida.csv", index=True)
```

/home/hugob/githu/Projecto2---Inteligencia-Artificial/venv/lib/python3.9/site-packages/pandas/core/frame.py:49
06: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
return super().drop()

Con Ridge Cross Validation (RidgeCV), podemos automatizar el proceso de encontrar el mejor alfa, en este caso lo buscamos con rangos entre 10⁻² y 10⁻³. Al utilizar el alfa que optimiza corremos la predicción con los datos de test, pero al subirlo a Kaggle, el puntaje no mejora.

```
In [384]: alpha_range = 10.**np.arange(-2, 3)
ridge = make_pipeline(StandardScaler(with_mean=False), RidgeCV(alpha=alpha_range, scoring='neg_mean_squared_error'))
ridge.fit(X_train, y_train)
print("*****Estadísticas training RidgeCV*****")
y_predicted = ridge.predict(X_train)
estadisticas(y_predicted, y_train)

print("*****Estadísticas testeo RidgeCV*****")
y_predicted = ridge.predict(X_holdout)
estadisticas(y_predicted, y_holdout)

*****Estadísticas training RidgeCV*****
R-Square Value 0.941438555977072

mean_absolute_error : 0.06750905576562512

mean_squared_error : 0.009207437100499215

root_mean_squared_error : 0.09595539120080339
*****Estadísticas testeo RidgeCV*****
R-Square Value 0.9077834704750518

mean_absolute_error : 0.08455365962958113

mean_squared_error : 0.01552133885940213

root_mean_squared_error : 0.12458466542637633
```

```
In [385]: test = X_dummies.loc["test"]
test.drop("SalePrice", axis=1, inplace=True)
y_predicted = ridge.predict(test)
y_predicted = np.exp(y_predicted)
Id = pd.read_csv("../Datos/test.csv").Id
salida = pd.DataFrame(y_predicted, Id, Id)
print(salida)
salida.to_csv("../Datos/salida_RidgeCV.csv", index=True)
```

```
0
1461 124117.840737
1462 118786.923146
1463 181346.566044
1464 201528.276367
1465 191117.112345
...
2915 84687.549599
2916 80734.269344
2917 169613.050229
2918 110683.687406
2919 226770.975831

[1459 rows x 1 columns]
```

/home/hugob/githu/Projecto2---Inteligencia-Artificial/venv/lib/python3.9/site-packages/pandas/core/frame.py:49
06: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
return super().drop()

Finalmente, probamos con regularización de Lasso, utilizando LassoCV para encontrar el mejor alfa y lo subimos a Kaggle. El puntaje en este caso es de 0.0001776256891830261. Corremos el modelo con este parámetro y subimos los datos al Kaggle. El puntaje RSME mejora, aunque no significativamente, el mismo se ubica en 0.14165.

salida_LassoCV.csv	0.14165
21 minutes ago by Hugo García	
Lasso	

```
In [386]: lassoregcv = LassoCV(n_alphas=100, normalize=True, random_state=1)
lassoregcv.fit(X_train, y_train)
print("alpha :", lassoregcv.alpha_)

print("*****Estadísticas training LassoCV*****")
y_predicted = lassoregcv.predict(X_train)
estadisticas(y_predicted, y_train)

print("*****Estadísticas testeo LassoCV*****")
y_predicted = lassoregcv.predict(X_holdout)
estadisticas(y_predicted, y_holdout)

*****Estadísticas training LassoCV*****
R-Square Value 0.926224369738492

mean_absolute_error : 0.07555574219992669

mean_squared_error : 0.011598773691243675

root_mean_squared_error : 0.10769760299674118
*****Estadísticas testeo LassoCV*****
R-Square Value 0.9164912230235194

mean_absolute_error : 0.0798919936495359

mean_squared_error : 0.014055701638994462

root_mean_squared_error : 0.11855674438425873
```

/home/hugob/githu/Projecto2---Inteligencia-Artificial/venv/lib/python3.9/site-packages/sklearn/linear_model/_base.py:141: FutureWarning: 'normalize' was deprecated in version 1.0 and will be removed in 1.2. If you wish to scale the data, use Pipeline with a StandardScaler in a preprocessing stage. To reproduce the previous behavior:

from sklearn.pipeline import make_pipeline
model = make_pipeline(StandardScaler(with_mean=False), Lasso())

If you wish to pass a sample_weight parameter, you need to pass it as a fit parameter to each step of the pipeline as follows:

kwargs = {'sample_weight': sample_weight for s in model.steps}
model.fit(X, y, **kwargs)

Set parameter alpha to: original_alpha * np.sqrt(n_samples).

warnings.warn()

```
In [387]: test = X_dummies.loc["test"]
test.drop("SalePrice", axis=1, inplace=True)
y_predicted = lassoregcv.predict(test)
y_predicted = np.exp(y_predicted)
Id = pd.read_csv("../Datos/test.csv").Id
salida = pd.DataFrame(y_predicted, Id, Id)
print(salida)
salida.to_csv("../Datos/salida_LassoCV.csv", index=True)
```

```
0
1461 123072.210011
1462 154171.888269
1463 176272.221306
1464 193615.612475
1465 196731.871988
...
2915 87445.728624
2916 88112.488746
2917 169632.770636
2918 118378.782318
2919 234475.091769

[1459 rows x 1 columns]
```

/home/hugob/githu/Projecto2---Inteligencia-Artificial/venv/lib/python3.9/site-packages/pandas/core/frame.py:49
06: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
return super().drop()

5. Conclusiones

1. Es importante hacer una buena exploración de los datos antes de comenzar con los modelos. Conocer los datos permitirá comprender de mejor manera el problema y abordarlo con las herramientas adecuadas.
2. La regresión lineal múltiple utilizando todos los parámetros produce un sobreentrenamiento, arrojando un modelo que es peor que si solo utilizáramos el promedio como predictor.
3. Al utilizar regularización de Ridge, se logra reducir la función de error, lo que permite tener un modelo más flexible. También es importante conocer la herramienta RidgeCV para encontrar de forma automatizada mejores parámetros de penalización.
4. La regularización de Lasso tiene la bondad de "eliminar" algunas variables que no son significativas en el modelo (pone coeficientes a cero). En esta parte sería importante analizar que variables están quitando para ver si concuerdan con la intuición construida en la parte exploratoria.
5. Es importante continuar mejorando el modelo y en su momento utilizar las técnicas que originalmente proponía el ejercicio, estos son random forest y gradient boosting.

6. Referencias

1. Gerón, Aurélien. Hands-on Machine Learning with Scikit-Learn, Keras & TensorFlow. O'Reilly Media. 2019.
2. Gareth, James; Witten, Daniela; Hastie, Trevor y Tibshirani, Robert. An Introduction to Statistical Learning. Second Edition.
3. <https://towardsdatascience.com/super-simple-machine-learning-by-me-multiple-linear-regression-part-1-447800e8b624>
4. <https://towardsdatascience.com/introduction-to-data-preprocessing-in-machine-learning-a9fa3a5dc9d>
5. <https://www.youtube.com/watch?v=q-DyJA8ZmYM>
6. <https://github.com/valaguer/machine-learning-mil2>
7. <https://harish-reddy.medium.com/regularization-in-python-699c3bad8622>