🔔

🏛 Getting Started Prediction Competition

# House Prices: Advanced Regression Techniques

Predict sales prices and practice feature engineering, RFs, and gradient boosting

**k** Kaggle · 5,142 teams · Ongoing

Overview    Data    Notebooks    Discussion    Leaderboard    Rules    Team                     My Submiss

## Overview

Description

Evaluation

Tutorials

Frequently Asked Questions

## Start here if...

You have some experience with R or Python and machine learning basics. This is a perfect competition for data science students who have completed an online course in machine learning and are looking to expand their skill set before trying a featured competition.

## Competition Description



Ask a home buyer to describe their dream house, and they probably won't begin with the height of the basement ceiling or the proximity to an east-west railroad. But this playground competition's dataset proves that much more influences price negotiations than the number of bedrooms or a white-picket fence.

With 79 explanatory variables describing (almost) every aspect of residential homes in Ames, Iowa, this competition challenges you to predict the final price of each home.

## Practice Skills

- Creative feature engineering

# Data fields

Here's a brief version of what you'll find in the data description file.

- SalePrice - the property's sale price in dollars. This is the target variable that you're trying to predict.
- MSSubClass: The building class
- MSZoning: The general zoning classification
- LotFrontage: Linear feet of street connected to property
- LotArea: Lot size in square feet
- Street: Type of road access
- Alley: Type of alley access
- LotShape: General shape of property
- LandContour: Flatness of the property
- Utilities: Type of utilities available
- LotConfig: Lot configuration
- LandSlope: Slope of property
- Neighborhood: Physical locations within Ames city limits
- Condition1: Proximity to main road or railroad
- Condition2: Proximity to main road or railroad (if a second is present)
- BldgType: Type of dwelling
- HouseStyle: Style of dwelling
- OverallQual: Overall material and finish quality
- OverallCond: Overall condition rating
- YearBuilt: Original construction date
- YearRemodAdd: Remodel date
- RoofStyle: Type of roof
- RoofMatl: Roof material
- Exterior1st: Exterior covering on house
- Exterior2nd: Exterior covering on house (if more than one material)
- MasVnrType: Masonry veneer type
- MasVnrArea: Masonry veneer area in square feet
- ExterQual: Exterior material quality
- ExterCond: Present condition of the material on the exterior
- Foundation: Type of foundation
- BsmtQual: Height of the basement
- BsmtCond: General condition of the basement
- BsmtExposure: Walkout or garden level basement walls
- BsmtFinType1: Quality of basement finished area
- BsmtFinSF1: Type 1 finished square feet
- BsmtFinType2: Quality of second finished area (if present)
- BsmtFinSF2: Type 2 finished square feet
- BsmtUnfSF: Unfinished square feet of basement area
- TotalBsmtSF: Total square feet of basement area
- Heating: Type of heating
- HeatingQC: Heating quality and condition
- CentralAir: Central air conditioning
- Electrical: Electrical system
- 1stFlrSF: First Floor square feet
- 2ndFlrSF: Second floor square feet
- LowQualFinSF: Low quality finished square feet (all floors)
- GrLivArea: Above grade (ground) living area square feet

- BsmtFullBath: Basement full bathrooms
- BsmtHalfBath: Basement half bathrooms
- FullBath: Full bathrooms above grade
- HalfBath: Half baths above grade
- Bedroom: Number of bedrooms above basement level
- Kitchen: Number of kitchens
- KitchenQual: Kitchen quality
- TotRmsAbvGrd: Total rooms above grade (does not include bathrooms)
- Functional: Home functionality rating
- Fireplaces: Number of fireplaces
- FireplaceQu: Fireplace quality
- GarageType: Garage location
- GarageYrBlt: Year garage was built
- GarageFinish: Interior finish of the garage
- GarageCars: Size of garage in car capacity
- GarageArea: Size of garage in square feet
- GarageQual: Garage quality
- GarageCond: Garage condition
- PavedDrive: Paved driveway
- WoodDeckSF: Wood deck area in square feet
- OpenPorchSF: Open porch area in square feet
- EnclosedPorch: Enclosed porch area in square feet
- 3SsnPorch: Three season porch area in square feet
- ScreenPorch: Screen porch area in square feet
- PoolArea: Pool area in square feet
- PoolQC: Pool quality
- Fence: Fence quality
- MiscFeature: Miscellaneous feature not covered in other categories
- MiscVal: $Value of miscellaneous feature
- MoSold: Month Sold
- YrSold: Year Sold
- SaleType: Type of sale
- SaleCondition: Condition of sale

Introduction

This kernel is an attempt to use every trick in the books to unleash the full power of Linear Regression, including a lot of preprocessing and a look at several Regularization algorithms.

At the time of writing, it achieves a score of about 0.121 on the public LB, just using regression, no RF, no xgboost, no ensembling etc. All comments/corrections are more than welcome.

**[1]:**
```python
# Imports
import pandas as pd
import numpy as np
from sklearn.model_selection import cross_val_score, train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression, RidgeCV, LassoCV, ElasticN
from sklearn.metrics import mean_squared_error, make_scorer
from scipy.stats import skew
from IPython.display import display
import matplotlib.pyplot as plt
import seaborn as sns

# Definitions
pd.set_option('display.float_format', lambda x: '%.3f' % x)
%matplotlib inline
#njobs = 4
```

**[2]:**
```python
# Get data
train = pd.read_csv("../input/train.csv")
print("train : " + str(train.shape))
```

```
train : (1460, 81)
```

**[3]:**
```python
# Check for duplicates
idsUnique = len(set(train.Id))
idsTotal = train.shape[0]
idsDupli = idsTotal - idsUnique
print("There are " + str(idsDupli) + " duplicate IDs for " + str(idsTotal) +

# Drop Id column
train.drop("Id", axis = 1, inplace = True)
```
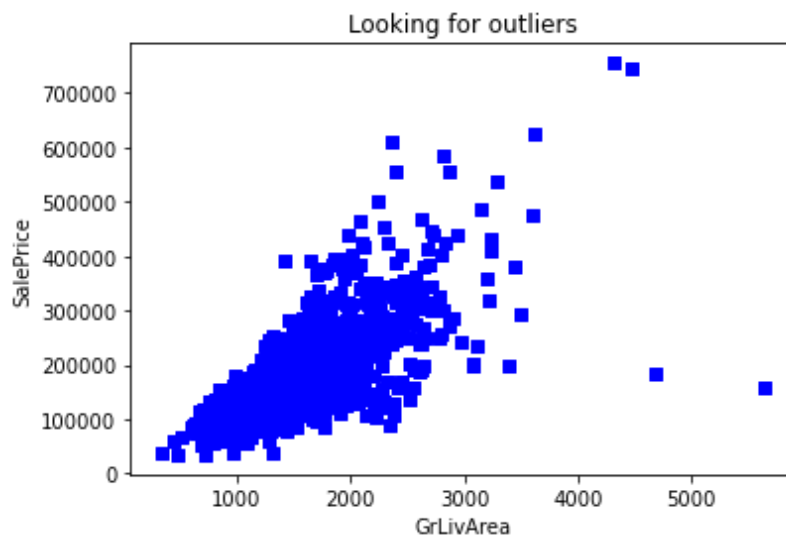
```
There are 0 duplicate IDs for 1460 total entries
```

Preprocessing

[4]:
```
# Looking for outliers, as indicated in https://ww2.amstat.org/publications/:
plt.scatter(train.GrLivArea, train.SalePrice, c = "blue", marker = "s")
plt.title("Looking for outliers")
plt.xlabel("GrLivArea")
plt.ylabel("SalePrice")
plt.show()

train = train[train.GrLivArea < 4000]
```



There seems to be 2 extreme outliers on the bottom right, really large houses that sold for really cheap. More generally, the author of the dataset recommends removing 'any houses with more than 4000 square feet' from the dataset.
Reference : https://ww2.amstat.org/publications/jse/v19n3/decock.pdf

[5]:
```
# Log transform the target for official scoring
train.SalePrice = np.log1p(train.SalePrice)
y = train.SalePrice
```

Taking logs means that errors in predicting expensive houses and cheap houses will affect the result equally.

[6]:

```python
# Handle missing values for features where median/mean or most common value c

# Alley : data description says NA means "no alley access"
train.loc[:, "Alley"] = train.loc[:, "Alley"].fillna("None")
# BedroomAbvGr : NA most likely means 0
train.loc[:, "BedroomAbvGr"] = train.loc[:, "BedroomAbvGr"].fillna(0)
# BsmtQual etc : data description says NA for basement features is "no baseme
train.loc[:, "BsmtQual"] = train.loc[:, "BsmtQual"].fillna("No")
train.loc[:, "BsmtCond"] = train.loc[:, "BsmtCond"].fillna("No")
train.loc[:, "BsmtExposure"] = train.loc[:, "BsmtExposure"].fillna("No")
train.loc[:, "BsmtFinType1"] = train.loc[:, "BsmtFinType1"].fillna("No")
train.loc[:, "BsmtFinType2"] = train.loc[:, "BsmtFinType2"].fillna("No")
train.loc[:, "BsmtFullBath"] = train.loc[:, "BsmtFullBath"].fillna(0)
train.loc[:, "BsmtHalfBath"] = train.loc[:, "BsmtHalfBath"].fillna(0)
train.loc[:, "BsmtUnfSF"] = train.loc[:, "BsmtUnfSF"].fillna(0)
# CentralAir : NA most likely means No
train.loc[:, "CentralAir"] = train.loc[:, "CentralAir"].fillna("N")
# Condition : NA most likely means Normal
train.loc[:, "Condition1"] = train.loc[:, "Condition1"].fillna("Norm")
train.loc[:, "Condition2"] = train.loc[:, "Condition2"].fillna("Norm")
# EnclosedPorch : NA most likely means no enclosed porch
train.loc[:, "EnclosedPorch"] = train.loc[:, "EnclosedPorch"].fillna(0)
# External stuff : NA most likely means average
train.loc[:, "ExterCond"] = train.loc[:, "ExterCond"].fillna("TA")
train.loc[:, "ExterQual"] = train.loc[:, "ExterQual"].fillna("TA")
# Fence : data description says NA means "no fence"
train.loc[:, "Fence"] = train.loc[:, "Fence"].fillna("No")
# FireplaceQu : data description says NA means "no fireplace"
train.loc[:, "FireplaceQu"] = train.loc[:, "FireplaceQu"].fillna("No")
train.loc[:, "Fireplaces"] = train.loc[:, "Fireplaces"].fillna(0)
# Functional : data description says NA means typical
train.loc[:, "Functional"] = train.loc[:, "Functional"].fillna("Typ")
# GarageType etc : data description says NA for garage features is "no garage
train.loc[:, "GarageType"] = train.loc[:, "GarageType"].fillna("No")
train.loc[:, "GarageFinish"] = train.loc[:, "GarageFinish"].fillna("No")
train.loc[:, "GarageQual"] = train.loc[:, "GarageQual"].fillna("No")
train.loc[:, "GarageCond"] = train.loc[:, "GarageCond"].fillna("No")
train.loc[:, "GarageArea"] = train.loc[:, "GarageArea"].fillna(0)
train.loc[:, "GarageCars"] = train.loc[:, "GarageCars"].fillna(0)
# HalfBath : NA most likely means no half baths above grade
train.loc[:, "HalfBath"] = train.loc[:, "HalfBath"].fillna(0)
# HeatingQC : NA most likely means typical
train.loc[:, "HeatingQC"] = train.loc[:, "HeatingQC"].fillna("TA")
# KitchenAbvGr : NA most likely means 0
train.loc[:, "KitchenAbvGr"] = train.loc[:, "KitchenAbvGr"].fillna(0)
# KitchenQual : NA most likely means typical
train.loc[:, "KitchenQual"] = train.loc[:, "KitchenQual"].fillna("TA")
# LotFrontage : NA most likely means no lot frontage
train.loc[:, "LotFrontage"] = train.loc[:, "LotFrontage"].fillna(0)
# LotShape : NA most likely means regular
train.loc[:, "LotShape"] = train.loc[:, "LotShape"].fillna("Reg")
# MasVnrType : NA most likely means no veneer
train.loc[:, "MasVnrType"] = train.loc[:, "MasVnrType"].fillna("None")
train.loc[:, "MasVnrArea"] = train.loc[:, "MasVnrArea"].fillna(0)
# MiscFeature : data description says NA means "no misc feature"
train.loc[:, "MiscFeature"] = train.loc[:, "MiscFeature"].fillna("No")
train.loc[:, "MiscVal"] = train.loc[:, "MiscVal"].fillna(0)
# OpenPorchSF : NA most likely means no open porch
train.loc[:, "OpenPorchSF"] = train.loc[:, "OpenPorchSF"].fillna(0)
# PavedDrive : NA most likely means not paved
train.loc[:, "PavedDrive"] = train.loc[:, "PavedDrive"].fillna("N")
```

```python
# PoolQC : data description says NA means "no pool"
train.loc[:, "PoolQC"] = train.loc[:, "PoolQC"].fillna("No")
train.loc[:, "PoolArea"] = train.loc[:, "PoolArea"].fillna(0)
# SaleCondition : NA most likely means normal sale
train.loc[:, "SaleCondition"] = train.loc[:, "SaleCondition"].fillna("Normal'
# ScreenPorch : NA most likely means no screen porch
train.loc[:, "ScreenPorch"] = train.loc[:, "ScreenPorch"].fillna(0)
# TotRmsAbvGrd : NA most likely means 0
train.loc[:, "TotRmsAbvGrd"] = train.loc[:, "TotRmsAbvGrd"].fillna(0)
# Utilities : NA most likely means all public utilities
train.loc[:, "Utilities"] = train.loc[:, "Utilities"].fillna("AllPub")
# WoodDeckSF : NA most likely means no wood deck
train.loc[:, "WoodDeckSF"] = train.loc[:, "WoodDeckSF"].fillna(0)
```

[7]:
```python
# Some numerical features are actually really categories
train = train.replace({"MSSubClass" : {20 : "SC20", 30 : "SC30", 40 : "SC40",
                                       50 : "SC50", 60 : "SC60", 70 : "SC70",
                                       80 : "SC80", 85 : "SC85", 90 : "SC90",
                                       150 : "SC150", 160 : "SC160", 180 : "S
                       "MoSold" : {1 : "Jan", 2 : "Feb", 3 : "Mar", 4 : "Apr"
                                   7 : "Jul", 8 : "Aug", 9 : "Sep", 10 : "Oct
                      })
```

```
[8]:    # Encode some categorical features as ordered numbers when there is informati
        train = train.replace({"Alley" : {"Grvl" : 1, "Pave" : 2},
                               "BsmtCond" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3,
                               "BsmtExposure" : {"No" : 0, "Mn" : 1, "Av": 2, "Gd" :
                               "BsmtFinType1" : {"No" : 0, "Unf" : 1, "LwQ": 2, "Rec'
                                                 "ALQ" : 5, "GLQ" : 6},
                               "BsmtFinType2" : {"No" : 0, "Unf" : 1, "LwQ": 2, "Rec'
                                                 "ALQ" : 5, "GLQ" : 6},
                               "BsmtQual" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA": 3, '
                               "ExterCond" : {"Po" : 1, "Fa" : 2, "TA": 3, "Gd": 4, '
                               "ExterQual" : {"Po" : 1, "Fa" : 2, "TA": 3, "Gd": 4, '
                               "FireplaceQu" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" :
                               "Functional" : {"Sal" : 1, "Sev" : 2, "Maj2" : 3, "Maj
                                               "Min2" : 6, "Min1" : 7, "Typ" : 8},
                               "GarageCond" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3
                               "GarageQual" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3
                               "HeatingQC" : {"Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4,
                               "KitchenQual" : {"Po" : 1, "Fa" : 2, "TA" : 3, "Gd" :
                               "LandSlope" : {"Sev" : 1, "Mod" : 2, "Gtl" : 3},
                               "LotShape" : {"IR3" : 1, "IR2" : 2, "IR1" : 3, "Reg" :
                               "PavedDrive" : {"N" : 0, "P" : 1, "Y" : 2},
                               "PoolQC" : {"No" : 0, "Fa" : 1, "TA" : 2, "Gd" : 3, "E
                               "Street" : {"Grvl" : 1, "Pave" : 2},
                               "Utilities" : {"ELO" : 1, "NoSeWa" : 2, "NoSewr" : 3,
                              )
```

Then we will create new features, in 3 ways :

1. Simplifications of existing features

2. Combinations of existing features

3. Polynomials on the top 10 existing features

[9]:

```python
# Create new features
# 1* Simplifications of existing features
train["SimplOverallQual"] = train.OverallQual.replace({1 : 1, 2 : 1, 3 : 1, #
                                                        4 : 2, 5 : 2, 6 : 2, #
                                                        7 : 3, 8 : 3, 9 : 3, 1
                                                        })
train["SimplOverallCond"] = train.OverallCond.replace({1 : 1, 2 : 1, 3 : 1, #
                                                        4 : 2, 5 : 2, 6 : 2, #
                                                        7 : 3, 8 : 3, 9 : 3, 1
                                                        })
train["SimplPoolQC"] = train.PoolQC.replace({1 : 1, 2 : 1, # average
                                             3 : 2, 4 : 2 # good
                                             })
train["SimplGarageCond"] = train.GarageCond.replace({1 : 1, # bad
                                                      2 : 1, 3 : 1, # average
                                                      4 : 2, 5 : 2 # good
                                                      })
train["SimplGarageQual"] = train.GarageQual.replace({1 : 1, # bad
                                                      2 : 1, 3 : 1, # average
                                                      4 : 2, 5 : 2 # good
                                                      })
train["SimplFireplaceQu"] = train.FireplaceQu.replace({1 : 1, # bad
                                                        2 : 1, 3 : 1, # averag
                                                        4 : 2, 5 : 2 # good
                                                        })
train["SimplFireplaceQu"] = train.FireplaceQu.replace({1 : 1, # bad
                                                        2 : 1, 3 : 1, # averag
                                                        4 : 2, 5 : 2 # good
                                                        })
train["SimplFunctional"] = train.Functional.replace({1 : 1, 2 : 1, # bad
                                                      3 : 2, 4 : 2, # major
                                                      5 : 3, 6 : 3, 7 : 3, # r
                                                      8 : 4 # typical
                                                      })
train["SimplKitchenQual"] = train.KitchenQual.replace({1 : 1, # bad
                                                        2 : 1, 3 : 1, # averag
                                                        4 : 2, 5 : 2 # good
                                                        })
train["SimplHeatingQC"] = train.HeatingQC.replace({1 : 1, # bad
                                                    2 : 1, 3 : 1, # average
                                                    4 : 2, 5 : 2 # good
                                                    })
train["SimplBsmtFinType1"] = train.BsmtFinType1.replace({1 : 1, # unfinished
                                                          2 : 1, 3 : 1, # rec
                                                          4 : 2, 5 : 2, 6 : 2
                                                          })
train["SimplBsmtFinType2"] = train.BsmtFinType2.replace({1 : 1, # unfinished
                                                          2 : 1, 3 : 1, # rec
                                                          4 : 2, 5 : 2, 6 : 2
                                                          })
train["SimplBsmtCond"] = train.BsmtCond.replace({1 : 1, # bad
                                                  2 : 1, 3 : 1, # average
                                                  4 : 2, 5 : 2 # good
                                                  })
train["SimplBsmtQual"] = train.BsmtQual.replace({1 : 1, # bad
                                                  2 : 1, 3 : 1, # average
                                                  4 : 2, 5 : 2 # good
                                                  })
train["SimplExterCond"] = train.ExterCond.replace({1 : 1, # bad
                                                    2 : 1, 3 : 1, # average
                                                    4 : 2, 5 : 2 # good
```

```
                                                      })
    train["SimplExterQual"] = train.ExterQual.replace({1 : 1, # bad
                                                        2 : 1, 3 : 1, # average
                                                        4 : 2, 5 : 2 # good
                                                      })


    # 2* Combinations of existing features
    # Overall quality of the house
    train["OverallGrade"] = train["OverallQual"] * train["OverallCond"]
    # Overall quality of the garage
    train["GarageGrade"] = train["GarageQual"] * train["GarageCond"]
    # Overall quality of the exterior
    train["ExterGrade"] = train["ExterQual"] * train["ExterCond"]
    # Overall kitchen score
    train["KitchenScore"] = train["KitchenAbvGr"] * train["KitchenQual"]
    # Overall fireplace score
    train["FireplaceScore"] = train["Fireplaces"] * train["FireplaceQu"]
    # Overall garage score
    train["GarageScore"] = train["GarageArea"] * train["GarageQual"]
    # Overall pool score
    train["PoolScore"] = train["PoolArea"] * train["PoolQC"]
    # Simplified overall quality of the house
    train["SimplOverallGrade"] = train["SimplOverallQual"] * train["SimplOverall(
    # Simplified overall quality of the exterior
    train["SimplExterGrade"] = train["SimplExterQual"] * train["SimplExterCond"]
    # Simplified overall pool score
    train["SimplPoolScore"] = train["PoolArea"] * train["SimplPoolQC"]
    # Simplified overall garage score
    train["SimplGarageScore"] = train["GarageArea"] * train["SimplGarageQual"]
    # Simplified overall fireplace score
    train["SimplFireplaceScore"] = train["Fireplaces"] * train["SimplFireplaceQu'
    # Simplified overall kitchen score
    train["SimplKitchenScore"] = train["KitchenAbvGr"] * train["SimplKitchenQual'
    # Total number of bathrooms
    train["TotalBath"] = train["BsmtFullBath"] + (0.5 * train["BsmtHalfBath"]) +
    train["FullBath"] + (0.5 * train["HalfBath"])
    # Total SF for house (incl. basement)
    train["AllSF"] = train["GrLivArea"] + train["TotalBsmtSF"]
    # Total SF for 1st + 2nd floors
    train["AllFlrsSF"] = train["1stFlrSF"] + train["2ndFlrSF"]
    # Total SF for porch
    train["AllPorchSF"] = train["OpenPorchSF"] + train["EnclosedPorch"] + \
    train["3SsnPorch"] + train["ScreenPorch"]
    # Has masonry veneer or not
    train["HasMasVnr"] = train.MasVnrType.replace({"BrkCmn" : 1, "BrkFace" : 1, '
                                                   "Stone" : 1, "None" : 0})
    # House completed before sale or not
    train["BoughtOffPlan"] = train.SaleCondition.replace({"Abnorml" : 0, "Alloca'
                                                          "Family" : 0, "Normal"
```

[10]:
```python
# Find most important features relative to target
print("Find most important features relative to target")
corr = train.corr()
corr.sort_values(["SalePrice"], ascending = False, inplace = True)
print(corr.SalePrice)
```

```
Find most important features relative to target
SalePrice          1.000
OverallQual        0.819
AllSF              0.817
AllFlrsSF          0.729
GrLivArea          0.719
                   ...
LandSlope         -0.040
SimplExterCond    -0.042
KitchenAbvGr      -0.148
EnclosedPorch     -0.149
LotShape          -0.286
Name: SalePrice, Length: 88, dtype: float64
```

[11]:
```python
# Create new features
# 3* Polynomials on the top 10 existing features
train["OverallQual-s2"] = train["OverallQual"] ** 2
train["OverallQual-s3"] = train["OverallQual"] ** 3
train["OverallQual-Sq"] = np.sqrt(train["OverallQual"])
train["AllSF-2"] = train["AllSF"] ** 2
train["AllSF-3"] = train["AllSF"] ** 3
train["AllSF-Sq"] = np.sqrt(train["AllSF"])
train["AllFlrsSF-2"] = train["AllFlrsSF"] ** 2
train["AllFlrsSF-3"] = train["AllFlrsSF"] ** 3
train["AllFlrsSF-Sq"] = np.sqrt(train["AllFlrsSF"])
train["GrLivArea-2"] = train["GrLivArea"] ** 2
train["GrLivArea-3"] = train["GrLivArea"] ** 3
train["GrLivArea-Sq"] = np.sqrt(train["GrLivArea"])
train["SimplOverallQual-s2"] = train["SimplOverallQual"] ** 2
train["SimplOverallQual-s3"] = train["SimplOverallQual"] ** 3
train["SimplOverallQual-Sq"] = np.sqrt(train["SimplOverallQual"])
train["ExterQual-2"] = train["ExterQual"] ** 2
train["ExterQual-3"] = train["ExterQual"] ** 3
train["ExterQual-Sq"] = np.sqrt(train["ExterQual"])
train["GarageCars-2"] = train["GarageCars"] ** 2
train["GarageCars-3"] = train["GarageCars"] ** 3
train["GarageCars-Sq"] = np.sqrt(train["GarageCars"])
train["TotalBath-2"] = train["TotalBath"] ** 2
train["TotalBath-3"] = train["TotalBath"] ** 3
train["TotalBath-Sq"] = np.sqrt(train["TotalBath"])
train["KitchenQual-2"] = train["KitchenQual"] ** 2
train["KitchenQual-3"] = train["KitchenQual"] ** 3
train["KitchenQual-Sq"] = np.sqrt(train["KitchenQual"])
train["GarageScore-2"] = train["GarageScore"] ** 2
train["GarageScore-3"] = train["GarageScore"] ** 3
train["GarageScore-Sq"] = np.sqrt(train["GarageScore"])
```

[12]:
```python
# Differentiate numerical features (minus the target) and categorical feature
categorical_features = train.select_dtypes(include = ["object"]).columns
numerical_features = train.select_dtypes(exclude = ["object"]).columns
numerical_features = numerical_features.drop("SalePrice")
print("Numerical features : " + str(len(numerical_features)))
print("Categorical features : " + str(len(categorical_features)))
train_num = train[numerical_features]
train_cat = train[categorical_features]
```

```
Numerical features : 117
Categorical features : 26
```

[13]:
```python
# Handle remaining missing values for numerical features by using median as
print("NAs for numerical features in train : " + str(train_num.isnull().value
train_num = train_num.fillna(train_num.median())
print("Remaining NAs for numerical features in train : " + str(train_num.isnu
```

```
NAs for numerical features in train : 81
Remaining NAs for numerical features in train : 0
```

[14]:
```python
# Log transform of the skewed numerical features to lessen impact of outliers
# Inspired by Alexandru Papiu's script : https://www.kaggle.com/apapiu/house-
# As a general rule of thumb, a skewness with an absolute value > 0.5 is cons
skewness = train_num.apply(lambda x: skew(x))
skewness = skewness[abs(skewness) > 0.5]
print(str(skewness.shape[0]) + " skewed numerical features to log transform")
skewed_features = skewness.index
train_num[skewed_features] = np.log1p(train_num[skewed_features])
```

```
86 skewed numerical features to log transform
```

[15]:
```python
# Create dummy features for categorical values via one-hot encoding
print("NAs for categorical features in train : " + str(train_cat.isnull().val
train_cat = pd.get_dummies(train_cat)
print("Remaining NAs for categorical features in train : " + str(train_cat.is
```

```
NAs for categorical features in train : 1
Remaining NAs for categorical features in train : 0
```

Modeling

[16]:
```python
# Join categorical and numerical features
train = pd.concat([train_num, train_cat], axis = 1)
print("New number of features : " + str(train.shape[1]))

# Partition the dataset in train + validation sets
X_train, X_test, y_train, y_test = train_test_split(train, y, test_size = 0.3
print("X_train : " + str(X_train.shape))
print("X_test : " + str(X_test.shape))
print("y_train : " + str(y_train.shape))
print("y_test : " + str(y_test.shape))
```

```
New number of features : 319
X_train : (1019, 319)
X_test : (437, 319)
y_train : (1019,)
y_test : (437,)
```

[17]:
```python
# Standardize numerical features
stdSc = StandardScaler()
X_train.loc[:, numerical_features] = stdSc.fit_transform(X_train.loc[:, numer
X_test.loc[:, numerical_features] = stdSc.transform(X_test.loc[:, numerical_f
```

```
/opt/conda/lib/python3.7/site-packages/pandas/core/indexing.py:1736: SettingWithCopyWarn
ing:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_
guide/indexing.html#returning-a-view-versus-a-copy
  isetter(loc, value[:, i].tolist())
/opt/conda/lib/python3.7/site-packages/pandas/core/indexing.py:1736: SettingWithCopyWarn
ing:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_
guide/indexing.html#returning-a-view-versus-a-copy
  isetter(loc, value[:, i].tolist())
```

Standardization cannot be done before the partitioning, as we don't want to fit the StandardScaler on
some observations that will later be used in the test set.

[18]:
```
# Define error measure for official scoring : RMSE
scorer = make_scorer(mean_squared_error, greater_is_better = False)

def rmse_cv_train(model):
    rmse= np.sqrt(-cross_val_score(model, X_train, y_train, scoring = scorer,
    return(rmse)

def rmse_cv_test(model):
    rmse= np.sqrt(-cross_val_score(model, X_test, y_test, scoring = scorer, c
    return(rmse)
```

Note : I'm not getting nearly the same numbers in local CV compared to public LB, so I'm a tad worried that my CV process may have an issue somewhere. If you spot something, please let me know.

1* Linear Regression without regularization

[19]:
```python
# Linear Regression
lr = LinearRegression()
lr.fit(X_train, y_train)

# Look at predictions on training and validation set
print("RMSE on Training set :", rmse_cv_train(lr).mean())
print("RMSE on Test set :", rmse_cv_test(lr).mean())
y_train_pred = lr.predict(X_train)
y_test_pred = lr.predict(X_test)

# Plot residuals
plt.scatter(y_train_pred, y_train_pred - y_train, c = "blue", marker = "s", l
plt.scatter(y_test_pred, y_test_pred - y_test, c = "lightgreen", marker = "s'
plt.title("Linear regression")
plt.xlabel("Predicted values")
plt.ylabel("Residuals")
plt.legend(loc = "upper left")
plt.hlines(y = 0, xmin = 10.5, xmax = 13.5, color = "red")
plt.show()

# Plot predictions
plt.scatter(y_train_pred, y_train, c = "blue", marker = "s", label = "Trainir
plt.scatter(y_test_pred, y_test, c = "lightgreen", marker = "s", label = "Val
plt.title("Linear regression")
plt.xlabel("Predicted values")
plt.ylabel("Real values")
plt.legend(loc = "upper left")
plt.plot([10.5, 13.5], [10.5, 13.5], c = "red")
plt.show()
```
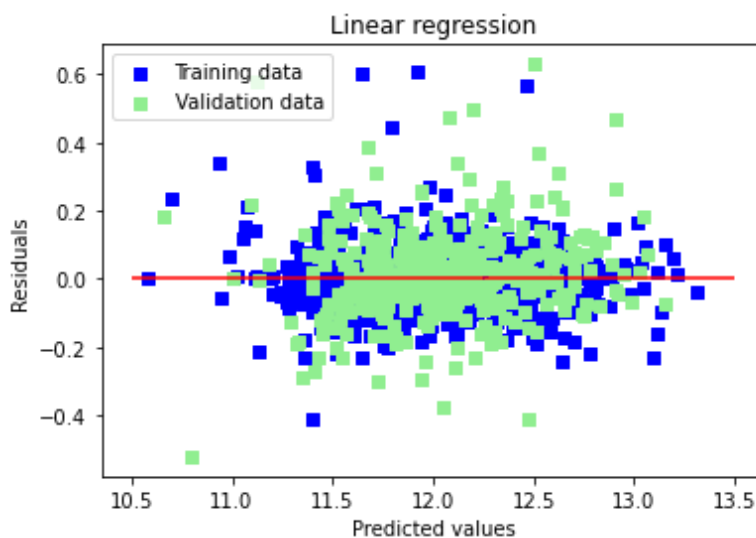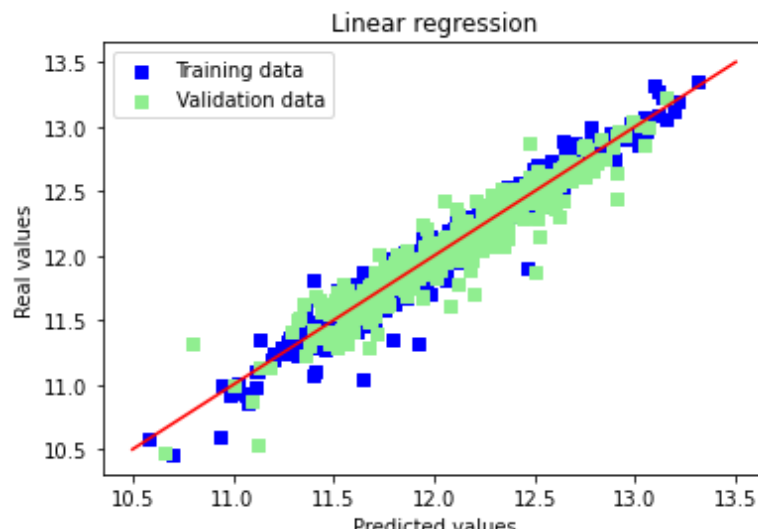
```
RMSE on Training set : 0.3885568698207793
RMSE on Test set : 0.40256620698831513
```

RMSE on Training set shows up weird here (not when I run it on my computer) for some reason. Errors seem randomly distributed and randomly scattered around the centerline, so there is that at least. It means our model was able to capture most of the explanatory information.