

Housing Prices

Ridge and Lasso regressions for the prediction of the median house price:
a machine learning implementation

Anna Olena Zhab'yak¹

University of Milan, Data science and Economics

and

Michele Maione²

University of Milan, Computer Science

(October 2020)

Key words: *machine learning, ridge regression, lasso regression, hedonic model*

Summary – This paper is trying to answer to the question what the best regression model is to predict the median house price. The hedonic theory is exploited and models as the Ridge and the Lasso regression are used on a cross sectional dataset of housing prices. Applying the cross-validation we estimate the parameters and then evaluate the model. Then PCA is implemented to improve the risk estimator. In section 1 we introduce the problem of housing prices in the U.S. and the approach used in this work. Section 2 is dedicated to the literature about regressions in predicting the price of real estates and 3rd section the theoretical notations are clarified to simplify the understanding of *the notions*. Our experiment is described in section 4 and in section 5 the consequential critical comments and evaluations.

1. Introduction and description of the problem

The hedonic theory identifies the attributes as implicitly embodied in goods and their observable market prices, so extending this concept to the housing prices we can see the attributes as the house's characteristics that are determinant for the final value. Hedonic model exploits the consumer theory and her willingness to pay depending on the utility gained from the bundle of aggregated attributes. Each attribute differently influences the price and its strength is given by the estimated coefficient.

Our work starts from a real problem of housing prices in the United States, where the economical purposes and the low mortgage rates incentive a solid and hot real estate market.³ Indeed, the U.S. is one of the most stable and secure countries for

¹ annaolena.zhab@studenti.unimi.it - 960298

² michele.maione@studenti.unimi.it - 931468

³ Santarelli, 2020

real estate investment in the recent years⁴. It is estimated that household wealth is nearly 50% invested in real estate and the owner-occupied housing rate in July 2019 was about 63.5%⁵. However, the U.S. real estate market was not always as reliable as today, indeed the sudden bubble of the housing market of 2006-2007 preceding the Great Recession and its subsequent burs is clear evidence of the system weaknesses. The speculation on the housing prices and their extremely high values is due to the lack of information caused by the manipulations of major players in the real estate sector⁶. For these reasons, the task of predicting the value of a house becomes crucial, as the constructed house price model can influence the economic growth and improve the efficiency of the real estate market. An accurate prediction model is significant and helps to fill up an information gap for the prospective homeowners, policy makers and other real estate market participants, such as, mortgage lenders and insurers⁷. Modelling house prices presents some issues, for example the median value might be extremely influenced by the value of the sold properties in the area with similar characteristics⁸ or the prediction could become wrong due to exogenous factors influencing the prices. Indeed, the economic health reflects in the market according to the supply and demand law so any shock will affect the current prices. Moreover, working on a large dataset, like the one used in this work, can lead to the so-called multicollinearity of the features which tend to overfit when it comes to implement the algorithm predicting the value. The classic OLS regression has the desired property of being unbiased, but it can suffer of overfitting and have a huge variance in those cases where features are highly correlated. To pull down the variance and obtain more biased estimator a regularization technique is necessary. The focus of

this paper is therefore on two regularization techniques, the Ridge and Lasso regression. The Ridge regression⁹ is a useful tool for improving prediction in regression tasks with highly correlated predictors. Lasso regression is also used to handle high dimensional databases where the features are correlated, and this technique shrinks some of them to zero, performing a feature selection with a consequent dimension reduction. Both methods act on the coefficients by introducing a penalty on them to make more effort to the most informative ones, this way minimizing overfitting of the data and solving the multicollinearity problem. The impact of each attribute on the predicted price is given by the value of the coefficient, higher coefficients mean higher influence. The penalty is the tool through which we perform the regularization, also called tuning parameter, it controls the bias-variance trade-off and the selection of it is crucial. For choosing the regularization parameter in practice, cross-validation (CV) is widely used.

2. Most important related works

Many works have been developed to predict the median house value with models of different complexity [see Manjula et al., 2017]. The concept of hedonic prices was developed by Rosen (1974), however the first the first implementing the hedonic model to the house sector was Lancaster (1966). Griliches (1971)¹⁰ provided the reading of a commodity, such as a house, as an aggregation of individual components or attributes. Timothy Oladunni & Sharad Sharma (2016) and Limsombunchai et al. (2004) have showed that the price of a property is predictable exploiting the hedonic theory, comparing the hedonic regression in comparison with other algorithms. Dubin (1998) has developed a work to predicted house prices using MLS data, even though exploiting different algorithms for the prediction, such as kriging algorithm

⁴ Source: International Investor Survey

⁵ Source: United States Census Bureau

⁶ Oladunni, Timothy & Sharma, Sharad, 2016

⁷ Limsombunchai et al. (2004)

⁸ the so-called Sales comparison approach

⁹ Introduced by Hoerl and Kennard (1970)

¹⁰ Griliches 1991

to create an accurate spatial interpolation of house prices. Others as Xin and Khalid (2018) have used ridge and lasso regression to deal with multicollinearity of features on a time series database for predicting the housing price. Hoerl and Kennard (1970) firstly introduced the Ridge regression as biased estimator for non-orthogonal problems. The asymptotic properties of ridge have been widely studied, [see for e.g. Dobriban and Wager (2018), Dicker (2016)]. For the validation approach we refer to the cross-validation which biased estimation of the error is known (Hastie et al., 2009, p. 243), since it uses a smaller amount of data than the entire dataset.¹¹ However, we can apply a bias-control, see Liu and Dobriban (2020), for example via k-fold cross validation, see Ray (2018), since there is an inverse relation between the k size and bias, if the first grows the latter goes down.

Mishra et al. (2017) have clearly explained the intuition behind the PCA and the underlying algebra to rich these results. PCA was introduced by Pearson (1901) and Hotelling (1933) and it is largely used in a lot of fields. Gupta and Kabundi (2010) have implemented lasso, ridge and PCA to predict housing prices on a time series dataset.

3. Notation and relevant definitions - Regression

The goal of the regression is to generate a prediction $\hat{y} = f(w, x)$ such that the loss function $\ell(y, \hat{y})$ is small for most data points $x \in \mathcal{X}$, where $\hat{y} \in \mathcal{Y}$ is the prediction from the labels set $\mathcal{Y} \subseteq \mathbb{R}$, $w \in \mathbb{R}^d$ is the coefficient vector and $\mathcal{X} = \mathbb{R}^d$ the data domain; the prediction mistakes are a function of the difference $|y - \hat{y}|$.

3.1. Hedonic model

Following the hedonic theory, the housing price can be written as a function $f(\cdot)$ in the following way:

$$P_i = f(s)$$

where s is the vector of all the objective attributes and P_i is the price of the i^{th} element of the data matrix \mathcal{X} . In this case, the price (our target variable) is a function of:

- longitude,
- latitude,
- housingMedianAge,
- totalRooms,
- totalBedrooms,
- population,
- households,
- medianIncome,
- medianHouseValue,
- oceanProximity.

3.2. Loss function

With loss function we denote the measure of how different the prediction of a hypothesis is from the true outcome. We use a nonnegative loss function to measure the discrepancy $\ell(y, \hat{y})$ between the predicted label \hat{y} and the true label y . In the regression task we define the quadratic loss that is the squared distance between y and \hat{y}

$$\ell(y, \hat{y}) = (y - \hat{y})^2$$

when $\hat{y} = y$ then $\ell(y, \hat{y}) = 0$ otherwise If $\hat{y} - y = c, \forall c \in \mathbb{R}^+$ and c is large then also $\ell(y, \hat{y})$ tend to be large. The mean of the squared error (MSE) will be used in the experiment.

3.3. Test error and training error

The split of dataset into two separate subsets is necessary in order to have some fresh data to estimate the predictive power of the algorithm. The dataset is divided in n elements for the test, and m elements for the training. Indeed, the validation is given by the test error which is:

$$\frac{1}{n} \sum_{t=1}^n \ell(y'_t, f(x'_t))$$

¹¹ In other words the algorithm has not enough data to train and be approximated

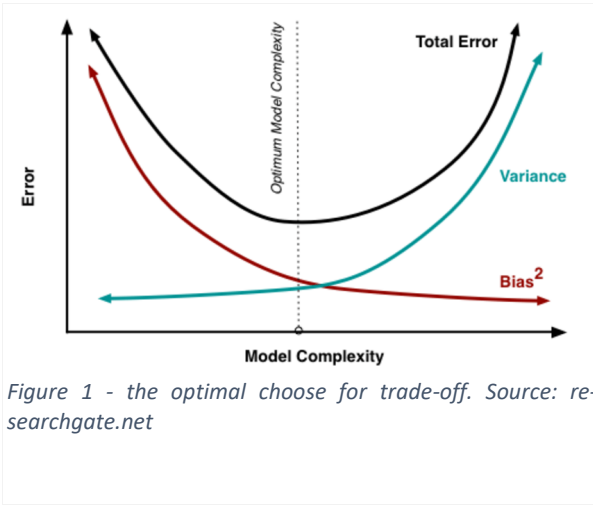
The validation is done over a fitted predictor in the training set, and its power is given by the training error:

$$\hat{\ell}(f) = \frac{1}{m} \sum_{t=1}^m \ell(y_t, f(x_t))$$

Total error is given by three elements:

- variance,
- bias,
- irreducible error.

The main idea is to derive a trade-off between the bias and variance, on order to optimize them both. More complex models present high variance and low bias since they fit good the true data but generalize worst.



3.4. Empirical Risk Minimization (ERM)

The Empirical Risk Minimization is a learning algorithm which returns some predictors $f \in \mathcal{F}$ given a set of predictors, that minimize the training error, given a non-negative real-valued loss function $\hat{\ell}$:

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\ell}(f)$$

3.5. Statistical risk, Bayes optimal predictor and Bayes optimal risk

We use statistical learning to introduce the notion of expectations in estimating the loss since we need to assume the independence between the variables and the predictor we generate is based

on this assumption. Let $h : \mathcal{X} \rightarrow \mathcal{Y}$ be the predictor that maps data points to labels. The **statistical risk** is then defined as the expectation of the loss function among D , the distribution from where the random sample of data points and labels were drawn:

$$\ell_D(h) = \mathbb{E}[\ell(Y, h(x))]$$

where $h(x)$ is the predicted \hat{y} . We then define as **Bayes optimal predictor** as the function f^* which minimize the overall training error $\ell_D(h)$, given the conditional probability among all predictors given that our data point is x :

$$f^*(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(Y, \hat{y}) | X = x]$$

The **Bayes optimal risk** is the expectation over the loss function of the Bayes optimal predictor and following the same logic as before we have that the Bayes risk is smaller than the other risks:

$$\mathbb{E}[\ell(Y, f^*(x))] \leq \mathbb{E}[\ell(Y, h(x))]$$

Coming to our regression problem with the squared loss, the Bayes optimal predictor is:

$$f^*(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[(Y - \hat{y})^2 | X = x] \quad (1)$$

minimizing this quantity¹², we have:

$$f^*(x) = \mathbb{E}[Y | X = x]$$

and the Bayes risk becomes the expectation of (1):

$$\mathbb{E}[(Y - f^*(X))^2 | X = x] = \operatorname{Var}[Y | X = x]$$

3.6. Regressions – Linear, Ridge, Lasso

Ridge and Lasso regression modify the standard linear regression by introducing a positive constant as regularization parameter. Indeed, the objective function to minimize under these solutions is **RSS¹³ + penalty**, and the penalty differs for the two methods. Starting from the classical linear model we have:

$$y_i = \bar{x}_i^T w$$

¹² By taking the derivative w.r.t \hat{y} , since $f^*(x)$ is differentiable

¹³ Sum of the squared residuals used for the classical OLS

let be the data domain $\mathcal{X} = \mathbb{R}^d$ and $x = (1, x_1, \dots, x_d)^{14}$ a row vector of \mathcal{X} . The linear predictor is a linear function $h : \mathbb{R}^d \rightarrow \mathbb{R}$, and for an activation function $f : \mathbb{R} \rightarrow \mathbb{R}$ we can write as follows:

$$h(x) = f(w^\top x)$$

where $w \in \mathbb{R}^d$ and $w^\top x = \sum_{i=1}^d w_i x_i$.

The Bayes optimal risk is given by

$$f^*(x) = \mathbb{E}[y | X = x]$$

and it is also an empirical risk minimization to $(x_1, y_1) \dots (x_m, y_m)$ is

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{t=1}^m (w^\top x_t - y_t)^2$$

Since we can rewrite these terms in vector notation, we have

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|v - y\|^2$$

for $v = (w^\top x_1, \dots, w^\top x_m)$ the vector of predictions and $y = (y_1, \dots, y_m)$ the vector of real labels and for $v, y \in \mathbb{R}^m$.

In matrix notation we have S the design matrix $S \in \mathbb{R}^{m \times d}$ with d features and m observations x_i that are rows of S^\top , and therefore the vector becomes $v = Sw$. Applying the ERM we derive

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|Sw - y\|^2$$

The solution to the ERM is the minimization of this convex function $F(w) = \|Sw - y\|^2$ using the **Euclidian distance**. To solve the problem in linear regression we can use the closed form solution, or the gradient descend.

If $S^\top S$ is a non-singular matrix¹⁵, and the conditions of the general position holds, the solution of the ERM is the closed form:

$$\nabla F(w) = 2S^\top(Sw - y) = 0$$

$$\hat{w} = (S^\top S)^{-1} S^\top y$$

In some cases, the linear regression performs well on the training data, having a low bias, but it gives a non-accurate estimate on different data. The reason why it occurs it is because of multicollinearity of the prediction vectors (as known as non-orthogonality)¹⁶. More in general with d large or n small, the risk that the model can overfit¹⁷ the data is high. The OLS estimator \hat{w} therefore is unbiased but have a huge variance and it is not stable. To overcome this problem, Ridge and Lasso regression help to prevent over-fitting which results from simple linear regression. We introduce a regularized parameter α which adds some bias¹⁸ whereas pushing the variance down. This also controls the model complexity, indeed the value of α has a direct relation with the complexity. This occurs to find the best trade-off between bias and variance to get to that sweet spot for having good predictive performance¹⁹.

The two methods work similarly but lead to different results, this happens because of the divergent formulas.

3.6.1. Ridge solution

Ridge regression uses the penalty multiplied by the square of the magnitude of the coefficients, also known as L2 regularization.

The ERM functional of Ridge regression is

$$\hat{w}_\alpha = \operatorname{argmin}_{w \in \mathbb{R}^d} \|Sw - y\|^2 + \alpha \|w\|^2 : \forall \alpha > 0$$

for $\alpha \rightarrow 0$, $\hat{w}_\alpha \rightarrow \hat{w}$ so the solution leads the linear regression, for $\alpha \rightarrow \infty$ the coefficient tend to a zero vector and the line becomes flatter, shrinking the linear regression solution towards to zero.

¹⁴ Add one extra feature to stabilize the prediction

¹⁵ This happens if the data points span $m \geq d$

¹⁶ Hoerl and Kennard, 2010

¹⁷ Overfitting: the algorithm performs very good on training data but cannot be generalized to a new bunch of data.

¹⁸ Bias is how well the fit correspond to the true value

¹⁹ See graphic 1 in this paper

To optimize the objective function, we take the gradient as before and solve for w to find a suitable value

$$\nabla F(w) = \|Sw - y\|^2 + \alpha \|w\|^2$$

$$2S^T Sw - S^T y + 2\alpha w = 0$$

$$(S^T S + \alpha I)w = S^T y^{20}$$

The new estimated parameter becomes:

$$\hat{w}_\alpha = (S^T S + \alpha I)^{-1} S^T y$$

This is the so called closed-form solution and α is the one measuring the stability of the procedure.

3.6.2. Lasso solution

Least Absolute Shrinkage and Selection Operator, or simply Lasso, is slightly different from the previous because the penalty is multiplied by the absolute value of the magnitude of coefficients, also known as L1 regularization

$$\hat{w}_{Lasso} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|Sw - y\|^2 + \alpha |w| : \forall \alpha > 0$$

For $\alpha \rightarrow \infty$, $\hat{w}_{Lasso} = 0$. The Lasso procedure encourages simple, sparse models²¹, indeed some coefficients can become zero and be eliminated from the model, this way performing a feature selection.

The shrinkage amount is given by the value of tuning parameter α . If α increase, we have some parameters go straightway to zero.

The optimization of a non-differentiable function as Lasso solution is done by a proximal gradient descend approach.

The first step is to take the gradient descend for current $w^{(k)}$ vector and form a new vector $z^{(k)}$:

$$z^{(k)} = w^{(k)} - \eta X^T (Xw^{(k)} - y)$$

Where η is the step size and k is the moment we are considering.

Then solve the proximal regularize problem for $w^{(k+1)}$ as follows:

$$w^{(k+1)} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|z^{(k)} - w\|_2^2 + \eta \alpha |w|$$

$$\alpha, \eta > 0$$

This is a scalar minimization problem indeed we can rewrite it as:

$$\min_{w_i : i=1, \dots, m} \sum_{i=1}^m \left(z_i^{(k)} - w_i \right)^2 + \alpha \eta |w_i|$$

Since we have an absolute value for $|w_i|$ we consider two cases:

Case 1 $w_i \geq 0$

$$\min_{w_i} \left(z_i^{(k)} - w_i \right)^2 + \alpha \eta w_i$$

differentiate with respect to w and solve:

$$-2(z_i - w_i)^2 + \partial \eta = 0$$

$$w_i = z_i - \frac{\alpha \eta}{2}$$

therefore, since we have the non-negativity constraint over w_i

$$\text{if } z_i > \frac{\alpha \eta}{2}, w_i = z_i - \frac{\alpha \eta}{2}$$

otherwise $w_i = 0$

Case 2 $w_i \leq 0$

$$\min_{w_i} \left(z_i^{(k)} - w_i \right)^2 - \alpha \eta w_i$$

differentiate with respect to w and solve for it:

$$-2(z_i - w_i)^2 - \partial \eta = 0$$

$$w_i = z_i + \frac{\alpha \eta}{2}$$

$$\text{if } z_i + \frac{\alpha \eta}{2} > 0, w_i = 0$$

$$\text{otherwise } w_i = z_i + \frac{\alpha \eta}{2}$$

The three solutions are also known as the "soft threshold" operation.

The update rule of the algorithm is therefore

The common point of these two methods is that adding the regularization parameter to the cost

²⁰ Adding the identity matrix fixes the invertibility problem, always compute inverse, and this is more stable solution

²¹ Stephanie Glen. "Lasso Regression: Simple Definition" From StatisticsHowTo.com: Elementary Statistics for the rest of us!

function the algorithm is forced to pick the lowest weights, indeed the goal is to ensure a small coefficient through this regularization parameter. The main difference is that many coefficients are exactly zeroed under lasso, which is never the case in ridge regression where there is not any elimination of coefficients. Moreover, Lasso arbitrarily selects any one feature among the highly correlated ones, leading to a higher variance than Ridge regression.

3.7. Cross-validation

Cross-validation (CV) is one of the techniques used to test the effectiveness of a machine learning models, it is also a re-sampling procedure used to evaluate a model if we have a limited data²². This approach however can be biased, therefore a K-fold cross-validation is largely used for evaluating the accuracy of model. This approach splits the dataset in two parts, the test and training sets. However, the training set is splatted into k-subsets and each fold in each interaction is used for testing while the remaining are used to training. This approach ensures that every observation from the original dataset has the chance of appearing in training and test set, this way decreasing the bias of the CV²³. This approach is used for both model selection and model validation. For model selection and therefore for hyperparameter tuning the error on the testing part of each fold is computed:

$$\hat{\ell}_{D_k}(h_k) = \frac{K}{m} \sum_{(x,y) \in D_k} \ell(y, h_k(x))$$

Where the number of subsets obtained is k, and D_k is the testing part, so out subsets of training part is D_1, \dots, D_k .

$$\min_{\theta \in \Theta} \ell_D(h_S^{(\theta)})$$

where $h_S^{(\theta)} = A_\theta(S)$, S is the training set, θ the tuning parameter, and A the learning algorithm.

For the model validation, given a fixed hyperparameter, we want to estimate the $\mathbb{E}[\ell_D(A)]$.

After k-fold cross validation, we will get k different model estimation errors, so the CV final error will be the

$$\frac{1}{K} \sum_{k=1}^K \hat{\ell}_{D_k}(h_k)$$

given the choice of two predictors, it repeatedly picks the most accurate of the two.

1. Shuffle the dataset randomly,
2. Split the dataset into k groups,
3. For each unique group:
 1. Take the group as a hold out or test data set,
 2. Take the remaining groups as a training data set,
 3. Fit a model on the training set and evaluate it on the test set,
 4. Retain the evaluation score and discard the model.
4. Summarize the skill of the model using the sample of model evaluation scores.

For tuning the parameter, nested cross validation is largely used. This kind of approach allows to have two loops in the CV, the inner loop is responsible for hyperparameter tuning while the outer loop is for error estimation.

3.8. Principal component analysis

Principal Component Analysis is a technique used for dimensionality reduction. Its goal is to reduce the number of features through a combination of the original data variables, in this way keeping most of the original information. Standardization

²³ As $k \rightarrow N$, leads to "Leave-one-out cross-validation", the leave one out approach, where the number of sets

equal the number of observations. On the contrary for smaller values of k we have the CV approach.

is needed. This technique finds the eigenvalues and eigenvectors of the correlation matrix²⁴.

The selection of the principal components is based on the variance caused to the target variable. The principal components will be then independent one from the other. The feature that cause more variance is the first principal component, and so on until we reach a suitable number of explained variance by the principal components.

A good way is to plot the variance against principal components and ignore the principal components with diminishing values as shown in the following graph:²⁵ This way we reduce the variance and we can improve the stability of the regression by solving the multicollinearity.

It is a way of identifying patterns in data and expressing the data in such a way as to highlight their similarities and differences.

4. Proof of a technical result

The demonstration of our experiment and all the critical considerations will be described here. This material is also available on GitHub at this url: <https://github.com/mikymaione/HousingPrices> we have created a Jupiter Notebook to illustrate the procedure followed step by step at this url: <https://github.com/mikymaione/HousingPrices/blob/master/Source%20code/HousingPrices/main.ipynb>

4.1. Data pre-processing

Before performing the analysis and regression, pre-processing of data is necessary. The dataset presents features that cannot be compared in a linear Euclidian space; therefore, geometry is not working properly on this row data. Indeed, to

learn the algorithm, we need to encode the features and raise them to a homogeneous level, so we can compare them²⁶. The dataset contains 20,640 observations and 10 features for each house including the median house value which is the target value that we are trying to predict. Firstly, we create the two constants of the target variable and designed matrix: X , y .

4.1.1. Missing values

The missing values must be handled to avoid errors in the execution of the code, so they are filled with the mean value of the corresponding column.

4.1.2. Categorical features

There is a categorical feature which represents the distance from the ocean, we transform the elements of the column into columns dummies and assign it to the data set²⁷. Even though among hedonic model literature²⁸ there is a concern about the statistical insignificance of some features such as household size, we use two approaches:

- keep all the features in regressions to have as much information as possible,
- use some unsupervised techniques to decide which feature to drop such as PCA.

4.1.3. Standardization

The standardization is done by subtracting the mean and divided by the variance. In this way we have $\mu = 0$ and $\sigma = 1$. This procedure is needed as both L1 and L2 assume that all features are centred around 0 and have unit variance. More in general this is important to those algorithms that use Euclidian distance and for PCA implementation.

²⁴ The eigenvectors of the Covariance matrix are actually the directions of the axes where there is the most variance (most information).

²⁵ Usman Malik

²⁶ For example, longitude and households are features with real numbers however they have different interpretations.

²⁷ This is called hot encoder technique; we do not worry about adding extra dimensions as the dummy variable sets to zero the features that do not belong to the given observation.

²⁸ Features such as garden, household size, neighborhood satisfaction and schools, as demonstrated by Berna and Craig.

4.1.4. Correlation matrix

We explore the correlation matrix between the features to see graphically how they move together. The used scoring is Pearson's coefficient r and $r \in [-1; 1]$, if the value is closer to 1 there is more correlation and the sign gives the direction of this correlation. Darker and lighter colours on this map are the two extremes.

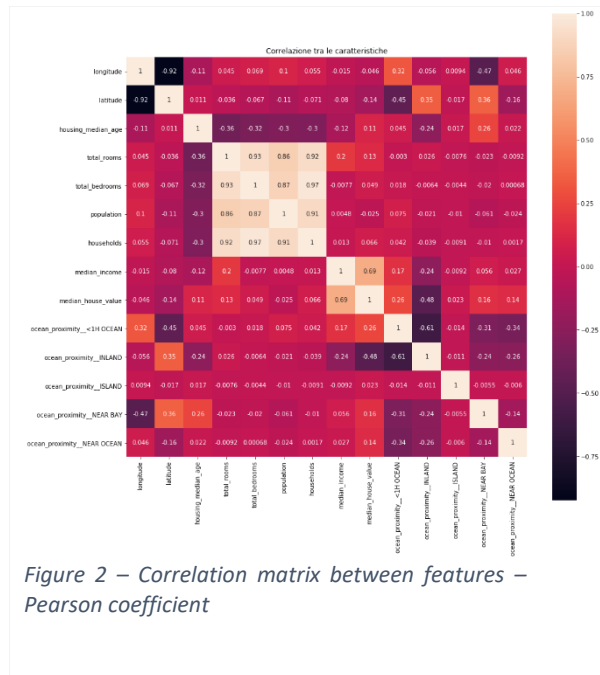


Figure 2 – Correlation matrix between features – Pearson coefficient

We have enough evidence that there exists statistical relationship between the variables²⁹. Therefore, our dataset is suitable for decomposition into its principal components to increase convergence speed and eliminate collinearity by finding the core components of the datasets.

4.2. Model tuning

Model tuning consists in the choice of the parameters to use into regression. In our case is the α hyperparameter for both Ridge and Lasso regression. The train-test split is done 80% - 20%³⁰

4.2.1. Scoring

Mean squared error is implemented as scoring, and it takes bigger values more than proportionally if the error in prediction increases.

4.2.2. Choosing the set of parameters alpha

In order to obtain a reasonable amount of information to determine a certain $f^* \in \mathcal{F}$ where f^* is the function that minimize the training error, we need to set different values of the tuning parameter to find out the best one. A larger value of α leads to a high bias but a low variance. On the other hand, for small values of α the variance increases, and bias go down. We perform the analysis on α with the relative mean squared error on the training data, comparing the Lasso and Ridge solutions.

Ridge regression

We use a logarithmic range $\ln|F|$. Training size m is bigger than the $\ln|F|$ in order to avoid underfitting.

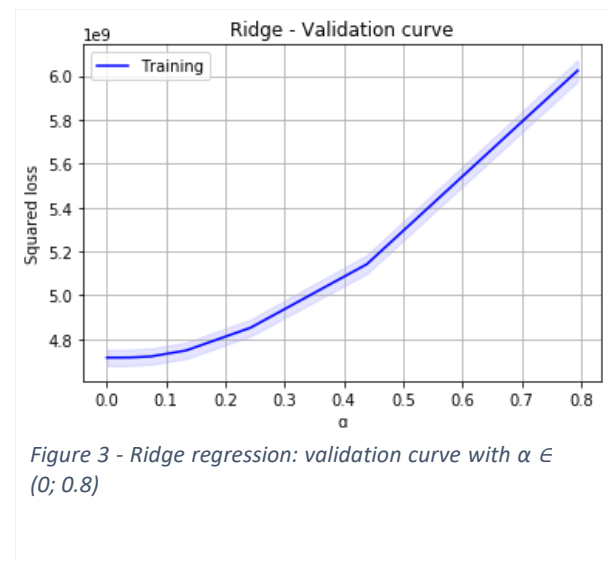


Figure 3 - Ridge regression: validation curve with $\alpha \in (0; 0.8)$

As we can see the optimal value for the hyper-parameter that optimize the squared loss is in between the range (0; 0.1], after that the squared lost increases. The best value of the penalized term is **0.000059**.

Lasso regression

For Lasso implementation we have these values:

²⁹ See the lighter square 4x4 in the middle.

³⁰ 80/20 rule: following the Pareto principle.

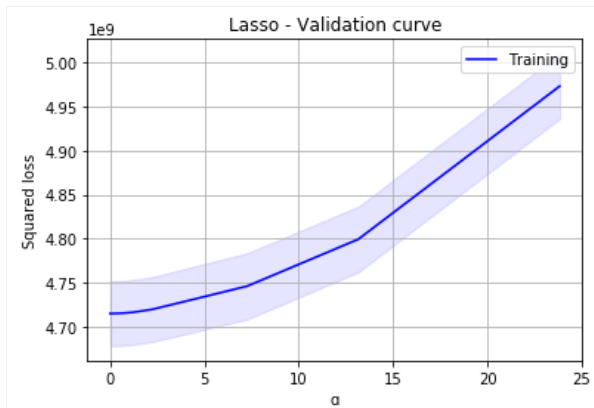


Figure 3 – Lasso regression: validation curve with $\alpha \in (0; 24)$

The values of alpha used are linear, and the best alpha according to the validation curve is **0.0003**.

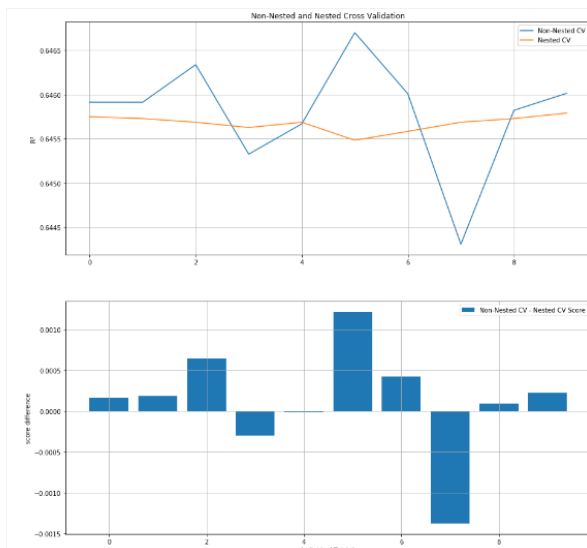


Figure 4 - LINEEEEEEE

4.3. Learning algorithm

Once we have determined the hyperparameter, the optimization of the learning algorithm is done. For ridge regression we use the Cholesky method, that is the closed form. For lasso regression we apply the proximate gradient descend³¹ This procedure is internal to the ridge and lasso function, and the learning algorithm is the output.

³¹ This is possible because we take an approximation of the gradient since Lasso function is not derivable.

4.3.1. Ridge learning algorithm

We fit the best α to plot the learning curve performance. Training error becomes larger when iterations are increased, and test error is higher as we could always expect a better performance on the training set. As we see the overfitting disappears as we increase the training size, around 5,000 training size, and it is improving with training size growth with a stable squared loss (4.7 in **Error! Reference source not found.**).

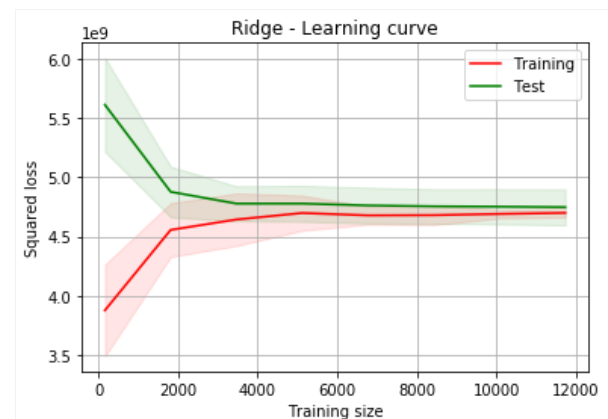


Figure 5 - Ridge regression: learning curve with different training set sizes

In this plot we can visualize how the predicted value differs from the real values. This is done with fit and predict functions. The prediction is more consistent with lower prices and becomes sparser for higher values, this can be caused by the presence of the outliers. The R^2 is around 63%.

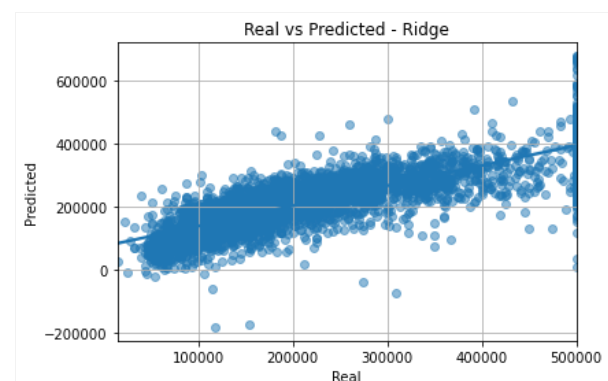


Figure 6 - Ridge regression: scatter plot prediction vs test

Here we can see the magnitude of each coefficient, and its prediction power on the target variable. The overall picture is very similar to the ridge coefficient, island location is the driven feature of the housing price.

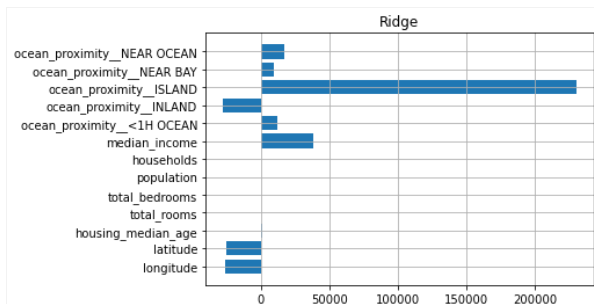


Figure 7 - Ridge: coefficients magnitude

4.3.2. Lasso learning algorithm

In lasso regression this increase is smoother, and MSE is more stable for α in between (0; 1]. This is what we could expect if there are feature highly correlated and that are not crucial for the regression. On the other hand, R^2 is specular to the MSE and tells us how well the model fits the data. The coefficient magnitude is shown in this graph. Median income has a huge power in the prediction, it means that it drives the values.

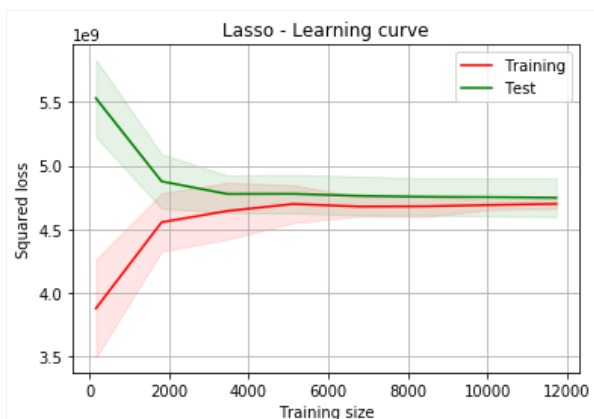


Figure 8 - Ridge regression: learning curve with different training set sizes

4.4. Principal Component Analysis

Cells that are highly correlated cluster together. Differences between the 1st pc (plotted on x axes)

are more important than the differences between the 2nd pc. Each point is a predictor that we have learned. Now we have a look on the Principal components of this data set and how is the variance distributed among two principal components.

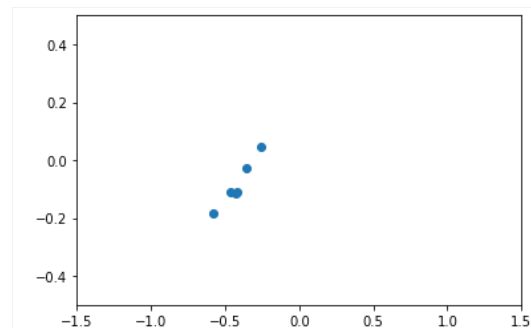


Figure 9 - PCA: Anna Olena scrivi qualcosa

After 5 features we do not gain more information therefore we will implement this decomposition.

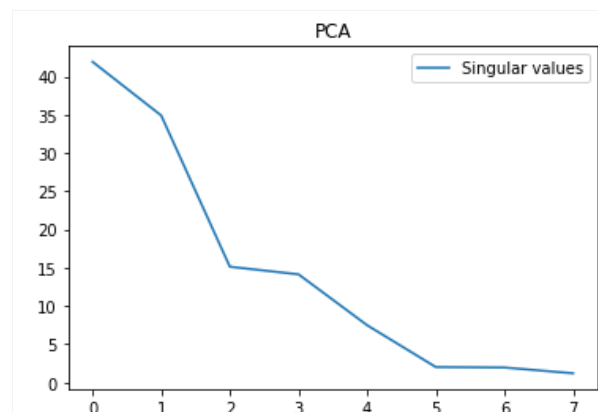


Figure 10 - PCA: singular values

This is the result and the performance of 5 PCA decomposition. The performance is not improving compared with the ridge regression.

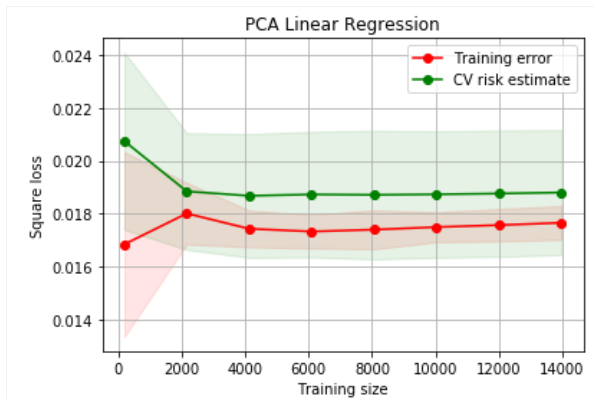


Figure 11 - PCA QUALCOSA

In the PCA analysis negative values of loadings of variable in the components of the PCA means the existence of an inverse correlation between the factor PCA and the variables.

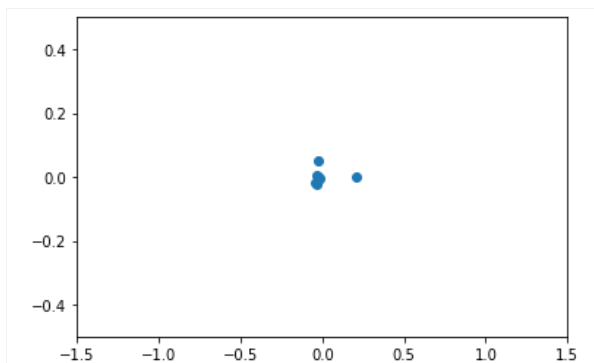


Figure 12 - PALLINI

5. Some critical considerations

As the model complexity increases, the models tends to fit even smaller deviations in the training data set. Though this leads to overfitting, let's keep this issue aside for some time and come to our main objective, i.e. the impact on the magnitude of coefficients. This can be analysed by looking at the data frame created above.

It is clear that the **size of coefficients increases exponentially with increase in model complexity**. I hope this gives some intuition into why putting a constraint on the magnitude of coefficients can be a good idea to reduce model complexity.

What does a large coefficient signify? It means that we are putting a lot of emphasis on that feature, i.e. the particular feature is a good predictor for the outcome. When it becomes too large, the algorithm starts modelling intricate relations to estimate the output and ends up overfitting to the training data.

Lasso can set some coefficients to zero, thus performing variable selection, while ridge regression cannot.

This way Lasso performs better in terms of reducing the variance in models with many redundant features.

Looking at the coefficients of ridge regression we can conclude that household size is statistically insignificant as the literature suggests. On the other hand, the driven force is the house location in island, which presents a direct positive impact on the predicted prices. Also the median income

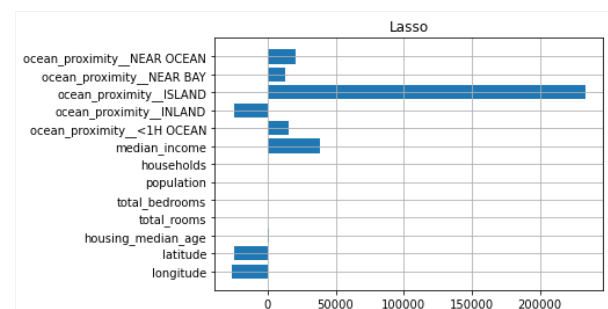


Figure 13 - Lasso: coefficients magnitude

In contrast, Ridge regression performs better in models where many features are important. This experiment shows that Lasso regression and PCA do not improve the risk estimate. This is because both reduce the dimension of the designed matrix. However, PCA performs better than Lasso, this is because the feature selection is done by keeping the most informative features, and Lasso regression just shrinks some coefficients to zero. In general, we can deduce that for this specific dataset the ridge regression is the most appropriate model, as if we drop some feature or elements the predictive power is poor indeed the performance

is better when all the features are considered in the prediction.

6. Bibliographical references

- Berna and Creig, 2016. "Defining spatial housing submarkets: Exploring the case for expert delineated boundaries". SAGE journals.
- Calhoun C. A., 2003, "Property Valuation Models and House Price Indexes for The Provinces of Thailand: 1992 –2000", *Housing Finance International*, 17(3): 31–41
- Dicker, Lee H. Ridge regression and asymptotic minimax estimation over spheres of growing dimension. *Bernoulli* 22 (2016), no. 1, 1--37. <https://projecteuclid.org/euclid.bj/1443620842>
- Dobriban, Edgar & Wager, Stefan. (2015). High-Dimensional Asymptotics of Prediction: Ridge Regression and Classification. *The Annals of Statistics*. 46
- Dubin, Robin. (1998). Predicting House Prices Using Multiple Listings Data. *The Journal of Real Estate Finance and Economics*. 17, 35-59
- Frew J. and B. Wilson, 2000, "Estimation The Connection Between Location and Property Value", *Essay in Honor of James A.Graaskamp*, Boston, MA: Kluwer Academic Publishers
- Frew J. and B. Wilson, 2000, "Estimation The Connection Between Location and Property Value", *Essay in Honor of James A.Graaskamp*, Boston, MA: Kluwer Academic Publishers
- Gupta, Kabundi, 2010, "Forecasting Real U.S.House Prices: Principal Components Versus Bayesian Regressions". *International Business & Economics Research Journal*.
- Hotelling H. 1933. Analysis of a complex of statistical variables into principal components. *J Educ Psychol*. 25: 417-441.
- Limsombunchai, Visit & Gan, Christopher & Lee, Minsoo. (2004). House Price Prediction: Hedonic Price Model vs. Artificial Neural Network. *American Journal of Applied Sciences*
- Liu, Sifan & Dobriban, Edgar. (2020). Ridge Regression: Structure, Cross-Validation, and Sketching
- Manjula, R & Jain, Shubham & Srivastava, Sharad & Kher, Pranav. (2017). Real estate value prediction using multivariate regression models. *IOP Conference Series: Materials Science and Engineering*. 263. 042098
- Mishra, Sidharth & Sarkar, Uttam & Taraphder, Subhash & Datta, Sanjoy & Swain, Devi & Saikhom, Reshma & Panda, Sasmita & Laishram, Menalsh. (2017). Principal Component Analysis. *International Journal of Livestock Research*.
- Oladunni, Timothy & Sharma, Sharad. (2016). Hedonic Housing Theory – A Machine Learning Investigation
- Pearson K. (1901) On lines and planes of closest fit to systems of points in space, *Philosophical Magazine* 2(11):559-572.
- Ray, 2018 www.analyticsvidhya.com/blog/2018/05/improve-model-performance-cross-validation-in-python-r
- Ridge Regression: Biased Estimation for Nonorthogonal Problems Author(s): Arthur E. Hoerl and Robert W. Kennard Source: *Technometrics*, Vol. 12, No. 1 (Feb., 1970), pp. 55-67 Published by: American Statistical Association and American Society for Quality Stable URL:
- Rosen S., 1974, "Hedonic Prices and Implicit Markets: Product Differentiation in Pure Competition", *Journal of Political Economics*, 82: 34 – 55
- Santarelli, 2020. "US Housing Market Forecast 2020 & 2021: Crash or Boom?". *Norada Real Estate*.
- Xin, Seng & Khalid, Kamil. (2018). Modeling House Price Using Ridge Regression and Lasso Regression. *International Journal of Engineering & Technology*. 7. 498

- Zvi Griliches, 1991. "Hedonic Price Indexes and the Measurement of Capital and Productivity: Some Historical Reflections," NBER Chapters, in: Fifty Years of Economic Measurement: The Jubilee of the Conference on Research in Income and Wealth, pages 185-206, National Bureau of Economic Research, Inc .

7. Sitographical references

- afire.org
- builtin.com
- census.gov
- datacamp.com
- jstor.org
- noradarealestate.com
- psu.edu
- researchgate.net
- stackabuse.com
- statisticsshowto.com
- towardsdatascience.com

8. Copyright

We declare that this material, which We now submit for assessment, is entirely our own work and has not been taken from the work of others, save and to the extent that such work has been cited and acknowledged within the text of our work. We understand that plagiarism, collusion, and copying are grave and serious offences in the university and accept the penalties that would be imposed should I engage in plagiarism, collusion or copying. This assignment, or any part of it, has not been previously submitted by us or any other person for assessment on this or any other course of study.

8.1. MIT License

Copyright (c) 2020 Anna Olena Zhab'yak, Michele Maione

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

The software is provided "as is", without warranty of any kind, express or implied, including but not limited to the warranties of merchantability, fitness for a particular purpose and noninfringement. In no event shall the authors or copyright holders be liable for any claim, damages or other liability, whether in an action of contract, tort or otherwise, arising from, out of or in connection with the software or the use or other dealings in the software.