

Ridge regression for housing prices

Giovanni Laganà, matr. 928792
giovanni.lagana@studenti.unimi.it

F943X - Statistical Methods for Machine Learning
Università degli Studi di Milano

September 29, 2020

I declare that this material, which I now submit for assessment, is entirely my 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 my work.

I 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 me or any other person for assessment on this or any other course of study.

1 Introduction

The aim of this project is to take advantage of **ridge regression** in order to create a model capable of inferring housing prices.

In general, many software libraries have been developed and are available today to do this task but in this project the implementation has been carried out from scratch.

The first part of the document introduces the theoretical tools that have been used, secondly, the dataset is presented; afterwards, it is shown how it has been preprocessed, and how ridge regression has been implemented.

The second part of the report deals with the different techniques that have been used to tune the regularization parameter, like cross validation and nested cross validation.

Moreover, principal component analysis has been performed with the aim of trying to reduce the cross-validated risk estimate.

Finally, results are presented, including a comparison with existing libraries.

2 Ridge regression: a supervised learning method

Ridge regression is a special case of linear regression and belongs to the family of supervised learning methods. We are, so, dealing with a numerical set of labels attached to data from the past and we want to predict a value which is as close as possible to the label of future unseen data.

Training, test and validation set Let us assume that we have dataset S organised as follows:

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

where $x_i \in \mathbb{R}^d$ is the vector of features of i -th row and $y_i \in \mathbb{R}$ is the single label of i -th row in the dataset.

S is randomly partitioned into training, validation and test set, with the purpose of isolating training, tuning and final tests on unseen data.

Training set The training set is used to fit the model.

Validation set The validation set is used to evaluate the performances of the trained model on new data, it is typically used to tune parameters in order to get the best configuration for the model that minimizes the loss function.

Test set The test set aim is to test the predictor on unseen data such that it is possible to evaluate its capability of generalization.

It is really important to split these sets such that the predictor never sees test set, otherwise *data leakage* phenomenon occurs, which leads to overfitting. Typically, this happens because the predictor tries to automatically adapt itself on the data which it shouldn't know and for which it is trained to work on.

Loss function It is a function $\ell(\hat{y}, y)$ that compares a predicted label \hat{y} to the real one y . Because of this, $\ell(\hat{y}, y)$ gives a quantifiable value of how much the prediction is close to the correct value.

The loss function is, typically, used to compute the training and the test errors.

There are plenty of loss functions, for this project we used the square loss¹:

$$\ell(\hat{y}, y) = (\hat{y} - y)^2 \quad (1)$$

Statistical risk We assume that the dataset S is a random sample (X, Y) drawn from an unknown joint probability distribution $\mathcal{X} \times \mathcal{Y}$, we call this distribution D .

We have a learning task, described by (D, ℓ) where ℓ is the loss function that we use and h is a predictor such that $h : X \rightarrow Y$.

$$\ell_D(h) = E[\ell(Y, h(X))] \quad (2)$$

defined as the expected value of loss function, the statistical risk measures the performances of a predictor h with respect to a given learning task (D, ℓ) ; since D is unknown, we can't compute exactly this value but we can estimate it through the test error ([4], [6]).

Test error is a sample mean of loss function and, thanks to the law of large numbers, we know that, with m increasing, sample mean tends to expected value.

MSE We use mean squared error, an error function that, through square loss, compute training and test error, according to which set we use:

$$MSE(\underline{\hat{y}}, \underline{y}) = \frac{1}{m} \ell(\underline{\hat{y}}, \underline{y}) = \frac{1}{m} \|\underline{\hat{y}} - \underline{y}\|^2 = \frac{1}{m} \sum_{j=1}^m (\hat{y}_j - y_j)^2 \quad (3)$$

This is going to be our risk estimator.

¹vectorial form: $\ell(\underline{\hat{y}}, \underline{y}) = \|\underline{\hat{y}} - \underline{y}\|^2$ where $\underline{\hat{y}}, \underline{y}$ are vectors respectively containing m predictions and m real labels

Linear regression Linear regression is the case where predictors are linear functions $h : \mathbb{R}^d \rightarrow \mathbb{R}$ where d is the number of features in input. The goal of regression is to predict a quantitative value, useful when labels are not a fixed symbolic set of classes (categorization).

The prediction is carried out by learning a linear relationship between features and the label to predict:

$$\hat{y}_i = \underline{w} \cdot \underline{x}_i \quad (4)$$

where \hat{y}_i is the single prediction of i -th row, we can do it $\forall i = 1 \dots m$ and compute $\hat{\underline{y}}$, the vector of all predictions; the intercept is considered by adding a feature to the dataset with constant value 1.

ERM Empirical risk minimization is used during training to fit the model through a closed form that allows us to compute quickly the weight vector that minimize ℓ on training set:

$$\hat{\underline{w}} = \arg \min_{\underline{w} \in \mathbb{R}^d} \ell(\hat{\underline{y}}, \underline{y}) = \arg \min_{\underline{w} \in \mathbb{R}^d} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \arg \min_{\underline{w} \in \mathbb{R}^d} \sum_{i=1}^m (\underline{w} \cdot \underline{x}_i - y_i)^2 = \arg \min_{\underline{w} \in \mathbb{R}^d} \|X\underline{w} - \underline{y}\|^2 \quad (5)$$

where $X = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_m]$ and is also known as **design matrix**, \underline{w} is the vector of weights and \underline{y} is the vector of labels in the training set.

However, this formula shows some issues: since we want to minimize $\|X\underline{w} - \underline{y}\|^2$ we compute the gradient and set the equation equal to 0, we obtain:

$$2X^T(X\underline{w} - \underline{y}) = 0 \iff X^T X \underline{w} = X^T \underline{y} \quad (6)$$

If $X^T X$ is non singular (invertible) then

$$\hat{\underline{w}} = (X^T X)^{-1} X^T \underline{y} \quad (7)$$

but $X^T X$ is non singular iff $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_m$ span \mathbb{R}^d , so $X^T X$ isn't always invertible, in this case we can't take advantage of this formula but we need an alternative approach (e.g. gradient descent). Moreover, whether it is non singular, this vector can be unstable: it can change a lot when the dataset is perturbed (e.g. changing a few examples), hence, instability leads to a variance error increase and so overfitting.

Regularization with ridge regression In order to fix the instability issue of the vector, we introduce ridge regression ([3], [7]), which consists of adding a regularization term:

$$\hat{\underline{w}}_\alpha = \arg \min_{\underline{w} \in \mathbb{R}^d, \alpha > 0} \|X\underline{w} - \underline{y}\|^2 + \alpha \|\underline{w}\|^2 \quad (8)$$

so if we now minimize it, we obtain:

$$2(X^T X \underline{w} - X^T \underline{y}) + 2\alpha \underline{w} = 0 \iff (X^T X + \alpha I) \underline{w} = X^T \underline{y} \quad (9)$$

such that the update will look like:

$$\hat{\underline{w}}_\alpha = (X^T X + \alpha I)^{-1} X^T \underline{y} \quad (10)$$

- when $\alpha \rightarrow 0$ we have standard linear regression solution
- when $\alpha \rightarrow \infty$ the solution becomes the zero vector

We can conclude that α can be used to control the bias error of the algorithm and prevent overfitting. Finally, there is no more issue about the singularity of $X^T X$: if $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ are the eigenvalues of $X^T X$, the eigenvalues of $\alpha I + X^T X$ are simply $\alpha + \lambda_1 \geq \dots \geq \alpha + \lambda_d \geq 0$ and, $\alpha I + X^T X$ is invertible for all $\alpha > 0$ ([3]).

3 Dataset

The dataset `cal-housing` contains information about 20640 houses in USA and it is composed by the following attributes:

- `longitude`: a measure of how far west a house is; a higher value is farther west
- `latitude`: a measure of how far north a house is; a higher value is farther north
- `housing_median_age`: median age of a house within a block; a lower number is a newer building
- `total_rooms`: total number of rooms within a block
- `total_bedrooms`: total number of bedrooms within a block
- `population`: total number of people residing within a block
- `households`: total number of households, a group of people residing within a home unit, for a block
- `median_income`: median income for households within a block of houses (measured in tens of thousands of US Dollars)
- `median_house_value`: median house value for households within a block (measured in US Dollars)
- `ocean_proximity`: location of the house w.r.t ocean/sea

Since we want to infer housing prices, we will consider `median_house_value` as label.

Data visualization In Figure 1 we can see the head of the dataset showing the first 10 rows

	<code>longitude</code>	<code>latitude</code>	<code>housing_median_age</code>	<code>total_rooms</code>	<code>total_bedrooms</code>	<code>population</code>	<code>households</code>	<code>median_income</code>	<code>median_house_value</code>	<code>ocean_proximity</code>
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252	452600.0	NEAR BAY
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014	358500.0	NEAR BAY
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574	352100.0	NEAR BAY
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431	341300.0	NEAR BAY
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462	342200.0	NEAR BAY
5	-122.25	37.85	52.0	919.0	213.0	413.0	193.0	4.0368	269700.0	NEAR BAY
6	-122.25	37.84	52.0	2535.0	489.0	1094.0	514.0	3.6591	299200.0	NEAR BAY
7	-122.25	37.84	52.0	3104.0	687.0	1157.0	647.0	3.1200	241400.0	NEAR BAY
8	-122.26	37.84	42.0	2555.0	665.0	1206.0	595.0	2.0804	226700.0	NEAR BAY
9	-122.25	37.84	52.0	3549.0	707.0	1551.0	714.0	3.6912	261100.0	NEAR BAY

Figure 1: Dataset

The main benefit of this step is that we quickly realised the presence of a categorical column: `ocean_proximity`; recalling the definition of regression, we must provide to the model an input of numerical values (\mathbb{R}^d), therefore, we will handle this issue.

Outliers detection The graphic tool we chose in order to detect outliers at first sight is the boxplot: in its definition, boxplot represents outliers as dots all those values that are

- beyond `max_threshold` = 3^{rd} quartile + $1.5 \cdot$ interquartile distance
- below `min_threshold` = 1^{st} quartile - $1.5 \cdot$ interquartile distance

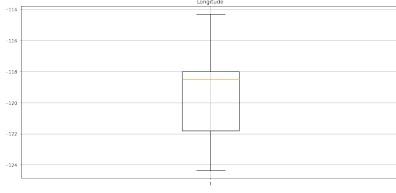


Figure 2: Longitude

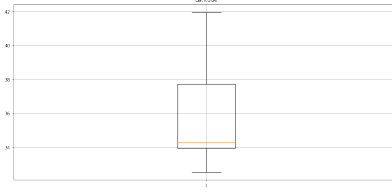


Figure 3: Latitude

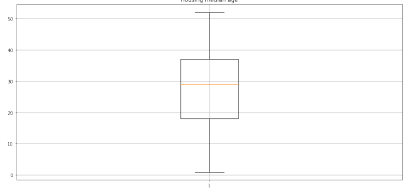


Figure 4: Median age

For the first three attributes (Figures 2, 3, 4), we don't see any outlier.

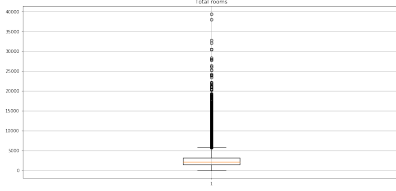


Figure 5: Total rooms

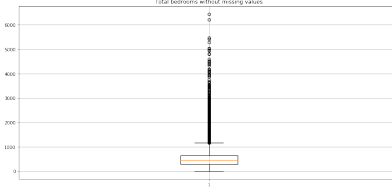


Figure 6: Total bedrooms

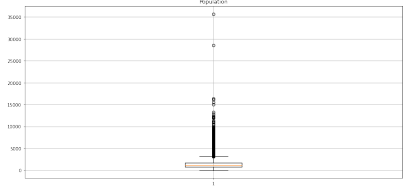


Figure 7: Population

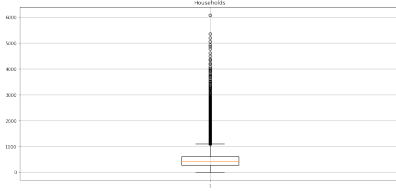


Figure 8: Households

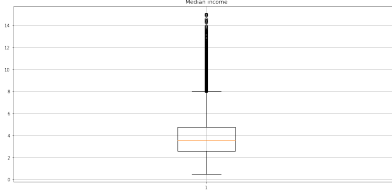


Figure 9: Median income

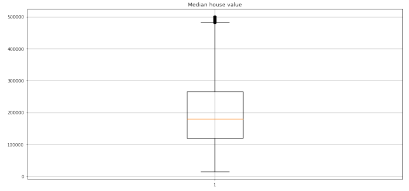


Figure 10: Median value

On the contrary, we notice outliers in the last six attributes, respectively `total_rooms`, `total_bedrooms`, `population`, `households`, `median_income` and `median_house_value` (Figures 5, 6, 7, 8, 9, 10).

4 Preprocessing

Once we detected anomalies, the next step is doing preprocessing.

Before applying the pipeline, we checked the presence of duplicates in the dataset, but luckily, there weren't any. The preprocessing pipeline is composed by:

- handling categorical values
- outliers removal
- handling missing values

- scaling
- adding intercept
- shuffling

Handling categorical values We decided to go ahead with one-hot encoding: a technique that replaces categorical values with truth values by adding as many columns as many categories and setting 1 if that value occurs, 0 otherwise.

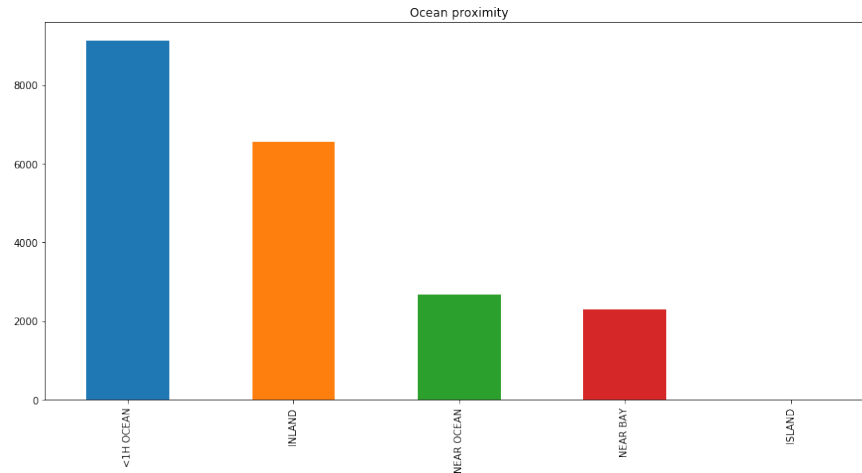


Figure 11: `ocean_proximity` distribution

We applied this encoding method for these reasons:

1. the number of possible values of `ocean_proximity` is small (see Figure 11)
2. an alternative label encoding forces us to choose discrete values according to some criterion (e.g. from the furthest to the closest w.r.t. the ocean) which could influence the results
3. the definition of the possible categories are in some cases vague and subjective for the dataset annotator (e.g. let us take `<1H OCEAN`, `NEAR OCEAN` and `NEAR BAY`: we don't know which category really indicates a closer location to the ocean)

the drawback is that linear models suffer immensely from the curse of dimensionality so adding features could be risky but, since we are talking about 4 more features (5 categories minus the drop of the original column) we think that this is the best trade off.

Outliers removal Outliers are values that show huge variance with respect to the rest of the distribution, this means that

- they could be due to an error
- they could be due to a legitimate property but they could mislead the regressor in catching rare dependencies between features and labels

in both cases leaving them wouldn't be an acceptable solution, since their presence may impact on the final behavior of our regressor, for this reason we decide to remove them. Outliers are contained in `total_rooms`, `total_bedrooms`, `population`, `households`, `median_income` and `median_house_value`; we had confirmed their presence through boxplots, we compute `max_threshold`, `min_threshold` and the `interquartile_distance` for those features so that we can filter outliers out.

<code>total_rooms</code>	<code>total_bedrooms</code>	<code>population</code>	<code>households</code>	<code>median_income</code>	<code>median_house_value</code>
6.23 %	6.16 %	5.79 %	5.91 %	3.30 %	5.19 %

Table 1: Percentage of outliers

We notice from Table 1 that the percentage of outliers varies from 3 to 6 percent with respect to the whole dataset, which is a non-negligible fraction.

Moreover, they are contemporary present in multiple columns at the same row: `median_house_value` and `median_income` for example, or the same row even includes missing values as well.

Totally, 3925 rows with at least one outlier have been filtered out, that is to say the 19.02 % of data².

Missing values Another anomaly which has been detected in this dataset is missing values, in particular they are contained only in `total_bedrooms` column, they actually represent the 0.97 % of the total dataset (a sample is taken in Figure 12).

	<code>longitude</code>	<code>latitude</code>	<code>housing_median_age</code>	<code>total_rooms</code>	<code>total_bedrooms</code>	<code>population</code>	<code>households</code>	<code>median_income</code>	<code>median_house_value</code>	<code><1H OCEAN</code>	<code>INLAND</code>
290	-122.16	37.77	47.0	1256.0	NaN	570.0	218.0	4.3750	161900.0	0	0
341	-122.17	37.75	38.0	992.0	NaN	732.0	259.0	1.6196	85100.0	0	0
563	-122.24	37.75	45.0	891.0	NaN	384.0	146.0	4.9489	247100.0	0	0
696	-122.10	37.69	41.0	746.0	NaN	387.0	161.0	3.9063	178400.0	0	0
738	-122.14	37.67	37.0	3342.0	NaN	1635.0	557.0	4.7933	186900.0	0	0
1097	-121.77	39.66	20.0	3759.0	NaN	1705.0	600.0	4.7120	158600.0	0	1
1456	-121.98	37.96	22.0	2987.0	NaN	1420.0	540.0	3.6500	204100.0	0	1
1493	-122.01	37.94	23.0	3741.0	NaN	1339.0	499.0	6.7061	322300.0	0	0
1606	-122.08	37.88	26.0	2947.0	NaN	825.0	626.0	2.9330	85000.0	0	0
2028	-119.75	36.71	38.0	1481.0	NaN	1543.0	372.0	1.4577	49800.0	0	1

Figure 12: Missing values

In order to handle missing values, we propose two solutions:

- a reduced dataset without missing values
- a dataset of same size but with replaced missing values

Our goal is to compare the final scores to see actually how much the two strategies differ.

Replacing missing values We used the stochastic regression imputation method: in general, `total_bedrooms` must be a subset of `total_rooms` so we can infer the value of the latter from the

²Note that this percentage is different than the sum of percentages in Table 1, which is 32.59 %, because many rows with multiple outliers and/or missing values have been found

former; the correlation coefficient between these two features is 0.89, confirming that there is a positive linear relationship.

Let us consider μ , σ respectively as mean and standard deviation of a distribution; we compute the distribution of fractions of total bedrooms with respect to the number of total rooms for the two vectors and then we compute their μ and σ . Our aim is to find a suitable range that describes how many bedrooms there are, on average, with respect to the number of rooms.

Since $\mu = 0.21$ and $\sigma = 0.05$, we can define a range $[\mu - \sigma, \mu + \sigma] = [0.16, 0.26]$ from which we can uniformly extract a number r_i at random $\forall i = 1 \dots g$ where g is the number of missing values; afterwards, we replace the i -th missing value with $total_rooms_i \cdot r_i$.

In order to keep consistency, we have to round this value to an integer value, since we are dealing with bedrooms and they belong to a discrete set³.

Scaling A general practice in machine learning is to scale features to a common range such that the application of some techniques can perform better (e.g. PCA), otherwise the different scale of each feature could impact on results. This approach is called scaling and, more specifically, normalization when the range is $[0, 1]$.

Moreover, scaling is usually done only on features and not on the target variable, this would allow us to keep trace of information about the error expressed in dollars. However, since there is a relevant difference in scale between the target variable and the rest of features, we decided to scale it as well, since it gets difficult to see how MSE changes every time⁴.

We used two different techniques: *min-max* and *zscore*:

Min-max Basically, it compresses values in the range $[0, 1]$.

$$\frac{value - min}{max - min}$$

We apply this formula for each row of each feature, for which we compute its minimum and maximum (Figure 13).

longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	<1H OCEAN	INLAND	ISLAND
0.213996	0.564293	1.000000	0.170526	0.089202	0.151095	0.150418	0.899633	0.721533	0.0	0.0	0
0.212982	0.564293	1.000000	0.132653	0.124413	0.171070	0.189415	0.684719	0.698417	0.0	0.0	0
0.212982	0.564293	1.000000	0.201923	0.159624	0.173325	0.226555	0.445496	0.700343	0.0	0.0	0
0.212982	0.564293	1.000000	0.062991	0.107199	0.124356	0.165274	0.470871	0.545164	0.0	0.0	0
0.212982	0.563231	1.000000	0.380102	0.323161	0.343750	0.463324	0.420587	0.608306	0.0	0.0	0
0.212982	0.563231	1.000000	0.491758	0.478091	0.364046	0.586815	0.348816	0.484590	0.0	0.0	0
0.211968	0.563231	0.803922	0.384027	0.460876	0.379832	0.538533	0.210414	0.453126	0.0	0.0	0
0.212982	0.563231	1.000000	0.579082	0.493740	0.490979	0.649025	0.424861	0.526756	0.0	0.0	0
0.211968	0.564293	1.000000	0.314757	0.280125	0.284472	0.359331	0.359880	0.570420	0.0	0.0	0
0.211968	0.564293	1.000000	0.570055	0.528951	0.475838	0.667595	0.368853	0.485446	0.0	0.0	0

Figure 13: After min-max normalization

³The order of operations in the pipeline is crucial: if we replace missing values and then remove outliers we risk to use outliers and create more of them; on the other hand, if we remove outliers and then replace missing values we have to be careful about the code type of missing values: in most of programming languages their presence influence operations of filtering like the one of outliers removal. Our solution was to temporary convert their type to an impossible integer (a negative number of bedrooms can't exists), safely remove outliers, and then replace missing values without outliers influence.

⁴We must recall that MSE is a relative measure that is scale-dependent: by scaling the target variable, we have closer MSEs that are easily comparable otherwise big scale pushes MSE to diverge to $+\infty$.

Zscore It centers the involved distribution in the origin of the Cartesian plane, let us denote μ , σ as mean and standard deviation of a fixed feature, the range will be approximately $[-3\sigma, +3\sigma]$:

$$\frac{value - \mu}{\sigma}$$

again, we apply it to each feature of the dataset (Figure 14).

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	<1H OCEAN	INLAND
2	-1.319182	0.999125	1.869305	-0.752144	-1.305628	-1.305431	-1.324942	2.547547	1.743141	-0.882688	-0.696125
3	-1.324179	0.999125	1.869305	-0.938804	-1.092339	-1.197896	-1.109645	1.419003	1.627918	-0.882688	-0.696125
4	-1.324179	0.999125	1.869305	-0.597401	-0.879051	-1.185755	-0.904601	0.162805	1.637520	-0.882688	-0.696125
5	-1.324179	0.999125	1.869305	-1.282142	-1.196614	-1.449390	-1.242924	0.296052	0.864035	-0.882688	-0.696125
6	-1.324179	0.994511	1.869305	0.280771	0.111556	-0.268235	0.402558	0.032005	1.178764	-0.882688	-0.696125
7	-1.324179	0.994511	1.869305	0.831078	1.050025	-0.158966	1.084330	-0.344875	0.562110	-0.882688	-0.696125
8	-1.329175	0.994511	1.044328	0.300114	0.945751	-0.073978	0.817773	-1.071651	0.405279	-0.882688	-0.696125
9	-1.324179	0.994511	1.869305	1.261459	1.144820	0.524404	1.427780	0.054446	0.772284	-0.882688	-0.696125
10	-1.329175	0.999125	1.869305	-0.041290	-0.149130	-0.587373	-0.171567	-0.286780	0.989927	-0.882688	-0.696125
11	-1.329175	0.999125	1.869305	1.216971	1.358108	0.442886	1.530302	-0.239662	0.566377	-0.882688	-0.696125

Figure 14: After zscore normalization

Adding intercept As mentioned before, we add a feature of constant value 1 to the datasets with the aim of including the intercept.

Shuffling Sometimes, datasets can show some order in the data, this obviously affects the results. Since the split of training and test set can be unlucky, training and test errors could be very different and wrongly interpreted.

We noticed some kind of geographical order in this dataset, it can be seen in `ocean_proximity`: rows tend to have some locality principle, hence, there are a lot of sequences of rows with the same value (see Figure 15); in order to fix this issue we shuffle the four datasets.

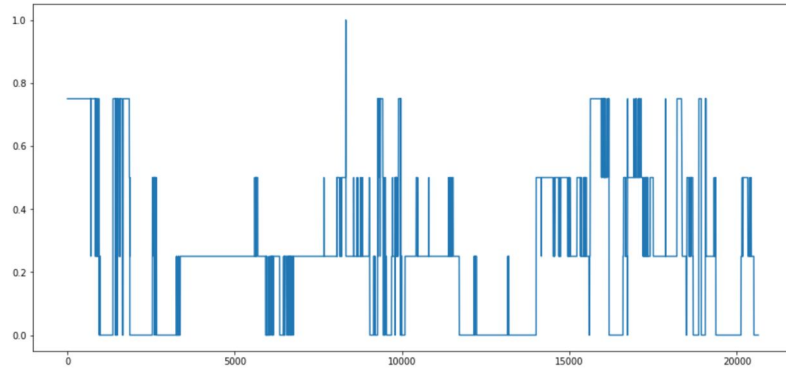


Figure 15: label-encoded `ocean_proximity` distribution

Any multicollinearity? We know that linear models suffer also from multicollinearity problem: the presence of multiple highly-correlated features could turn out to be redundant and also counterproductive.

For this reason, they are typically detected and removed, but this is not always the case ([2]).

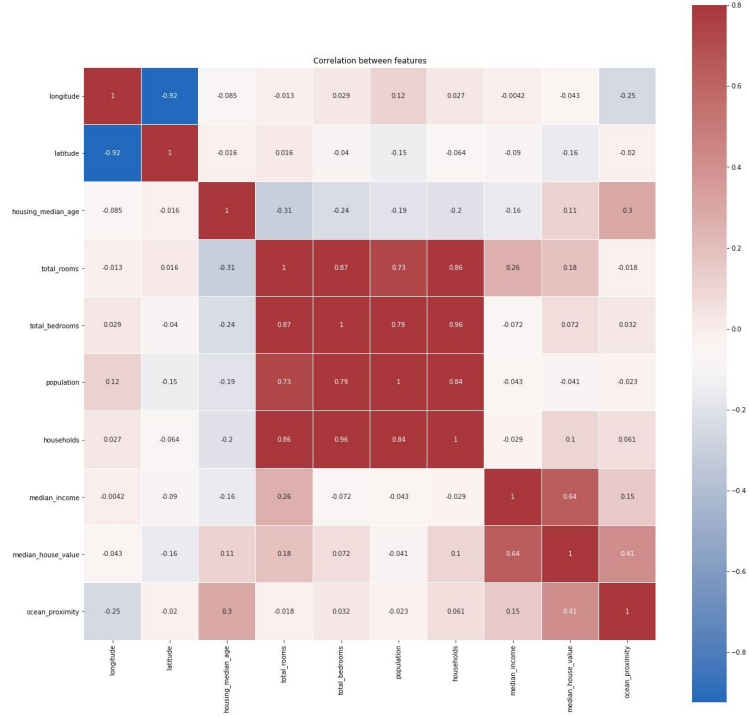


Figure 16: Correlation matrix between features

In Figure 16 we see that, actually, there are some highly correlated features but in this case it is better to keep them for three reasons:

1. results get worse with their removal
2. we used the correlation between `total_rooms`; and `total_bedrooms` to impute missing values
3. the regularization parameter should fix multicollinearity problems ([1])

At the end of preprocessing pipeline, we obtain 4 datasets:

1. `REP_MINMAX`: replaced missing values and min-max normalization
2. `RED_MINMAX`: reduced missing values and min-max normalization
3. `REP_ZSCORE`: replaced missing values and scaled with zscore
4. `RED_ZSCORE`: reduced missing values and scaled with zscore

5 Experiments

After preprocessing, the part of experiments is discussed in this section.

Model The model has been implemented with the following structure:

- a `fit` procedure (see Algorithm 1)
- a `predict` procedure (see Algorithm 2)

Algorithm 1 `fit` procedure

Procedure `fit`(X, \underline{y}, α)

Input: X is the design matrix of training set, \underline{y} is the vector of labels, α is the regularization parameter

Output: weight vector \underline{w}

- 1: $\underline{w} = (X^T X + \alpha I)^{-1} X^T \underline{y}$
 - 2: **return** \underline{w}
-

Algorithm 2 `predict` procedure

Procedure `predict`(X, \underline{w})

Input: X is the design matrix of test set, \underline{w} is the weight vector

Output: predictions vector $\hat{\underline{y}}$

- 1: $\hat{\underline{y}} = X \underline{w}$
 - 2: **return** $\hat{\underline{y}}$
-

Tuning hyperparameters Tuning hyperparameters means finding the best configuration for those variable whose value control the learning process, in this case α .

In order to find the best α of the model we started running multiple executions of the algorithm with different values of α (Algorithm 3).

Algorithm 3 training, validation and test error computation

- 1: training, validation, test = `split`(dataset)
 - 2: **for** α in $\underline{\alpha}$ **do**
 - 3: $\underline{w} = \text{fit}(\text{training.X}, \text{training.y}, \alpha)$
 - 4: $\hat{\underline{y}} = \text{predict}(\text{training.X}, \underline{w})$
 - 5: training_error = $\text{MSE}(\hat{\underline{y}}, \text{training.y})$
 - 6: $\hat{\underline{y}} = \text{predict}(\text{validation.X}, \underline{w})$
 - 7: validation_error = $\text{MSE}(\hat{\underline{y}}, \text{validation.y})$
 - 8: $\hat{\underline{y}} = \text{predict}(\text{test.X}, \underline{w})$
 - 9: test_error = $\text{MSE}(\hat{\underline{y}}, \text{test.y})$
 - 10: **end for**
-

At this stage, we randomly partitioned the four datasets in training, validation and test set with a

proportion of 60-20-20, whereas during the next paragraphs the cross-validation technique has been used, hence, no validation set was needed.

We created a vector of 5000 values from 0.1 to 500 and set it to α .

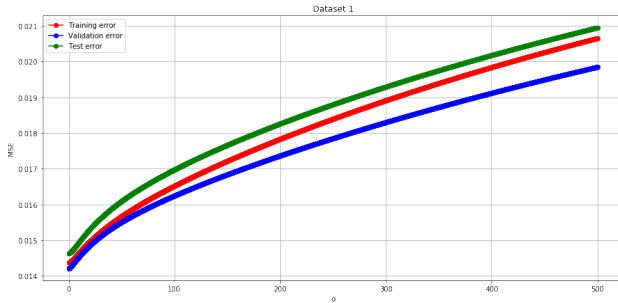


Figure 17: REP_MINMAX

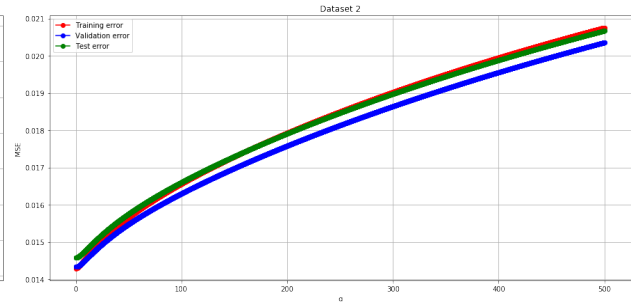


Figure 18: RED_MINMAX

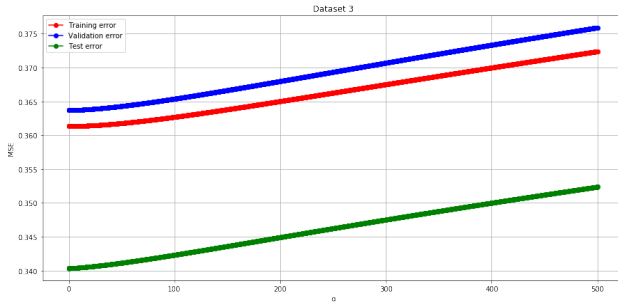


Figure 19: REP_ZSCORE

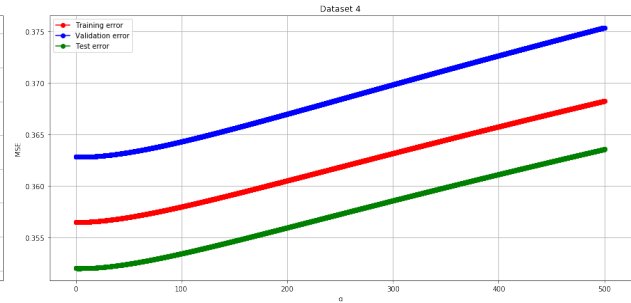


Figure 20: RED_ZSCORE

As we can see in Figures 17, 18, training error, as expected, tends to be lower than test error, since predictions are made on the set that the model has been trained on. However, in Figures 19, 20 it is shown the reversed case: test error is lower than train error, this suggests that it can happen that test set is composed by many more "easier" cases than the training part.

Regardless of this, we temporary ignore test error and look at training and validation errors, we see that with all four datasets the model perform better with α close to 0; this is the reason why we chose $\alpha = 0.1$ for most of the following experiments.

K-fold cross validation The rigorous definition tells us that when we refer to a learning algorithm $A(S)$ with S training set, we are implicitly meaning a family $A_\theta(S)$ of learning algorithms with $\theta \in \Theta$, the set of all possible hyperparameter values (in this case α).

Let us fix α and let ℓ_D be the statistical risk, we want to compute $E[\ell_D(A)]$ to understand the goodness of A .

We can do it through K -fold cross validation, a partitioning method that splits the dataset S in K folds S_1, S_2, \dots, S_K such that $\forall k = 1 \dots K$ we consider $S^{(k)} \equiv S \setminus S_k$ as training part and S_k as testing part.

Given that m is the size of S , K the number of folds, S_k the k -th fold and h_k the linear predictor returned by the training performed by A over S_k , then the scaled test error (in our case MSE) of each fold is

$$\hat{\ell}_{S_k}(h_k) = \frac{K}{m} \sum_{(x,y) \in S_k} \ell(h_k(x), y) \quad (11)$$

(11) estimates the statistical risk of the predictor h_k output by $A(S^{(k)})$ and it is done through MSE over S_k ; we compute the mean to estimate the expected risk of A

$$E[\ell_D(A)] = \frac{1}{K} \sum_{k=1}^K \hat{\ell}_{S_k}(h_k) \quad (12)$$

We did this in Algorithm 4, note that for a complete analysis, we computed both training and test errors:

Algorithm 4 K-fold cross validation procedure

Procedure $cv(S, K, \alpha)$

Input: S is the dataset, K the number of folds, α the regularization parameter

Output: training errors and test errors

```

1:  $S_1, S_2, \dots, S_K = \text{random\_partitioning}(S, K)$ 
2: for  $S_i$  in  $S_1, S_2, \dots, S_K$  do
3:    $\text{test} = S_i$ 
4:    $\text{train} = S \setminus S_i$ 
5:    $\underline{w} = \text{fit}(\text{train}.X, \text{train}.y, \alpha)$ 
6:    $\hat{y} = \text{predict}(\text{train}.X, \underline{w})$ 
7:    $\text{training\_error} = \text{MSE}(\hat{y}, \text{train}.y)$ 
8:    $\hat{y} = \text{predict}(\text{test}.X, \underline{w})$ 
9:    $\text{test\_error} = \text{MSE}(\hat{y}, \text{test}.y)$ 
10:   $\text{training\_errors.add}(\text{training\_error})$ 
11:   $\text{test\_errors.add}(\text{test\_error})$ 
12: end for
13: return  $\text{training\_errors}, \text{test\_errors}$ 
```

In general, K is chosen freely, but is typically recommended to be equal to 5 or 10.

Algorithm 5 cross-validated risk estimate computation

```

1:  $\hat{r}_1 = \text{mean}(cv(\text{REP\_MINMAX}, 5, 0.1).test\_errors)$ 
2:  $\hat{r}_2 = \text{mean}(cv(\text{RED\_MINMAX}, 5, 0.1).test\_errors)$ 
3:  $\hat{r}_3 = \text{mean}(cv(\text{REP\_ZSCORE}, 5, 0.1).test\_errors)$ 
4:  $\hat{r}_4 = \text{mean}(cv(\text{RED\_ZSCORE}, 5, 0.1).test\_errors)$ 
```

We report in Table 2 the cross-validated risk estimate with $K = 5$ and $\alpha = 0.1$ relative to the four datasets (Algorithm 5).

\hat{r}_1	\hat{r}_2	\hat{r}_3	\hat{r}_4
0.0144	0.0144	0.3583	0.3573

Table 2: Cross validation with $\alpha = 0.1$

Dependence of the cross-validated risk estimate on the parameter alpha With the aim of studying the dependence between them we use again the vector α containing 5000 values from 0.1 to 500 and we compute the cross-validated risk estimate for all α in α .

We do this for the four datasets, and we plot both training and test error:

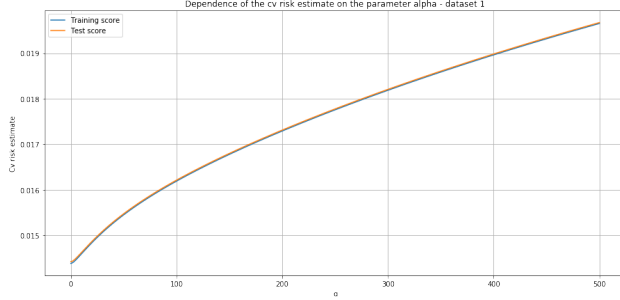


Figure 21: REP_MINMAX

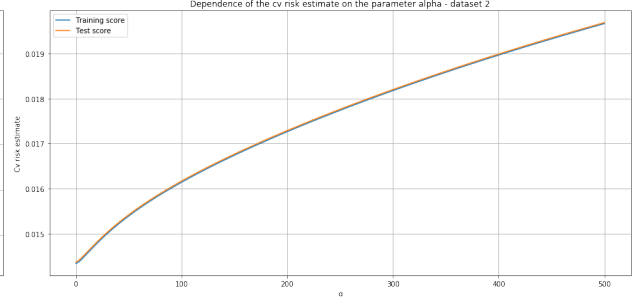


Figure 22: RED_MINMAX

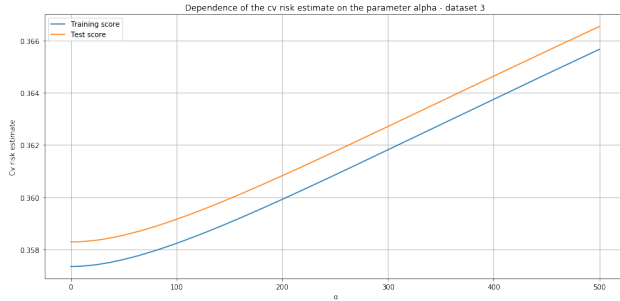


Figure 23: REP_ZSCORE

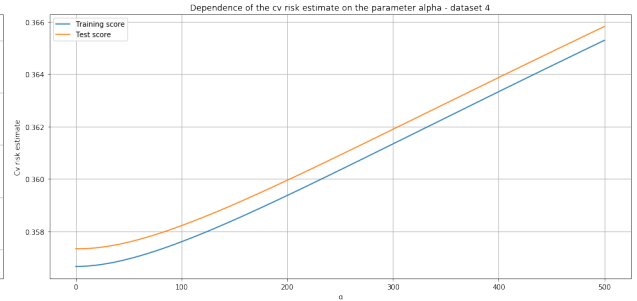


Figure 24: REP_ZSCORE

We clearly see that the cross-validated risk estimate increases with α and also decreases with α decreasing.

It means that there is a strong dependence of positive linear relationship between them, which is highly confirmed by a correlation coefficient of 0.99.

Concerning training and test errors, since the range of α is big, they are very close, especially in the dataset that has been normalized with min-max (see Figures 21, 22. On the other hand, the zscore dataset (Figures 23, 24) shows a better distinction of the two errors because the different range obtained with zscore allows MSE to be dilated.

Using PCA to reduce risk estimate Principal component analysis has been used with the goal of trying to reduce risk estimate.

PCA is a dimensionality-reduction method which is used to reduce the dimensionality of large datasets, by transforming a large set of variables into a smaller one that still contains most of the information in the large set.

Reducing the number of variables of a dataset naturally comes at the expense of accuracy, but the trick in dimensionality reduction is to trade a little accuracy for simplicity, because smaller datasets are easier to explore, visualize and make analyzing data much easier and faster for machine learning algorithms without extraneous variables to process ([5]).

We tried all possible reductions from 1 to 10 components but all of them didn't improve the results,

in Table 3 PCA decomposition with 5 components has been and cross-validated risk estimate is shown.

\hat{r}_1_{PCA}	\hat{r}_2_{PCA}	\hat{r}_3_{PCA}	\hat{r}_4_{PCA}
0.1556	0.1556	0.4568	0.4558

Table 3: Cross validation with PCA (same configuration of Table 2)

Nested cross validation It adds an inner layer of cross validation to the one we had previously defined in order to try a range of values for hyperparameter α and find the best one (Algorithm 6)

We define the vector $\underline{\alpha}$ containing 2500 values from 0.01 to 25, $K = 5$, $L = 4$; we use the procedure defined in Algorithm 6 to compute the nested cross-validated risk estimate (Algorithm 7)

we represent in Table 4 the final results.

$nest_r_1$	$\hat{\alpha}_1$	$nest_r_2$	$\hat{\alpha}_2$	$nest_r_3$	$\hat{\alpha}_3$	$nest_r_4$	$\hat{\alpha}_4$
0.0144	0.76	0.0144	0.01	0.3583	0.01	0.3573	0.01

Table 4: Nested cross validation

6 Results and comments

After the experiments some considerations can be done: the best performance is apparently achieved by the dataset normalized with min-max, though, as we said previously, MSE is a relative measure. For this reason, we must define how much low is a score to be considered "accurate": with min-max values have been shrunk between 0 and 1, this implies that loss can be at most 1.

On the other hand, zscore lets loss vary over a wider range, so in this case ℓ could reach, in the worst case, approximately 6σ , which is equal⁵ to 6.

Consequently, if we rescale the risk of REP_ZSCORE and RED_ZSCORE by dividing with 6, we get 0.05955 which is much more similar to the 0.0144 obtained by the first two datasets; anyway, we notice a little improvement using min-max.

Concerning the choice of removing or replacing missing values, there is no evident difference: probably, the impact of these values is so low because their fraction with respect to the whole dataset is small. In this context, only in the zscore dataset there is a slightly visible difference, but this is still due to the different scale that it has on MSE, in this case the reduced dataset performs a little better.

Finally, PCA didn't improve the risk estimate and nested cross validation confirmed that the best performances are achieved with α going towards 0.

⁵This value has been computed by calculating the mean of standard deviation of each feature, we obtained $\sigma = 1$, therefore $6\sigma = 6$

Algorithm 6 nested cross validation procedure

Procedure nest_cv($S, K, L, \underline{\alpha}$)

Input: S is the dataset, K folds, L subfolds, $\underline{\alpha}$ the vector of α values to try

Output: training errors, test errors and the best α

```
1:  $S_1, S_2, \dots, S_K = \text{partitioning}(S, K)$ 
2: for  $S_i$  in  $S_1, S_2, \dots, S_K$  do
3:   test =  $S_i$ 
4:   train =  $S \setminus S_i$ 
5:    $S_1, S_2, \dots, S_L = \text{partitioning}(\text{test}, L)$ 
6:   for  $S_j$  in  $S_1, S_2, \dots, S_L$  do
7:     in_test =  $S_j$ 
8:     in_train = train  $\setminus S_j$ 
9:     for  $\alpha$  in  $\underline{\alpha}$  do
10:       $\underline{w}^* = \text{fit}(\text{in\_train}.X, \text{in\_train}.\underline{y}, \alpha)$ 
11:       $\hat{\underline{y}}^* = \text{predict}(\text{in\_train}.X, \underline{w}^*)$ 
12:      in_training_error =  $\text{MSE}(\hat{\underline{y}}^*, \text{in\_train}.\underline{y})$ 
13:       $\hat{\underline{y}}^* = \text{predict}(\text{in\_test}.X, \underline{w}^*)$ 
14:      in_test_error =  $\text{MSE}(\hat{\underline{y}}^*, \text{test}.\underline{y})$ 
15:      in_training_errors.add(in_training_error)
16:      in_test_errors.add(in_test_error)
17:    end for
18:    best_in_result = min(in_test_errors)
19:     $\hat{\alpha}^* = \arg \min_{\alpha \in \underline{\alpha}} (\text{in\_test\_errors})$ 
20:    if best_in_result < best_result then
21:      best_result = best_in_result
22:       $\hat{\alpha} = \hat{\alpha}^*$ 
23:    end if
24:  end for
25:   $\underline{w} = \text{fit}(\text{train}.X, \text{train}.\underline{y}, \hat{\alpha})$ 
26:   $\hat{\underline{y}} = \text{predict}(\text{train}.X, \underline{w})$ 
27:  training_error =  $\text{MSE}(\hat{\underline{y}}, \text{train}.\underline{y})$ 
28:   $\hat{\underline{y}} = \text{predict}(\text{test}.X, \underline{w})$ 
29:  test_error =  $\text{MSE}(\hat{\underline{y}}, \text{test}.\underline{y})$ 
30:  training_errors.add(training_error)
31:  test_errors.add(test_error)
32: end for
33: return training_errors, test_errors,  $\hat{\alpha}$ 
```

Algorithm 7 nested cross-validated risk estimate and best alpha computation

```
1: result = nest_cv(score_REP_MINMAX, 5, 4,  $\alpha$ )
2:  $nest\_r_1^*$  = mean(result.test_errors)
3:  $\hat{\alpha}_1$  = result. $\hat{\alpha}$ 
4: result = nest_cv(score_RED_MINMAX, 5, 4,  $\alpha$ )
5:  $nest\_r_2^*$  = mean(result.test_errors)
6:  $\hat{\alpha}_2$  = result. $\hat{\alpha}$ 
7: result = nest_cv(score_REP_ZSCORE, 5, 4,  $\alpha$ )
8:  $nest\_r_3^*$  = mean(result.test_errors)
9:  $\hat{\alpha}_3$  = result. $\hat{\alpha}$ 
10: result = nest_cv(score_RED_ZSCORE, 5, 4,  $\alpha$ )
11:  $nest\_r_4^*$  = mean(result.test_errors)
12:  $\hat{\alpha}_4$  = result. $\hat{\alpha}$ 
```

7 External libraries

Scikit-learn is the classical example of software library which has developed many useful functions, including ridge regression implementation and cross validation.

We compute the cross-validated risk estimate again (Algorithm 8):

Algorithm 8 cross validation with *Scikit-learn* library

```
1:  $r_1^*$  = mean(sklearnCV(REP_MINMAX, sklearnRIDGE(0.1)))
2:  $r_2^*$  = mean(sklearnCV(RED_MINMAX, sklearnRIDGE(0.1)))
3:  $r_3^*$  = mean(sklearnCV(REP_ZSCORE, sklearnRIDGE(0.1)))
4:  $r_4^*$  = mean(sklearnCV(RED_ZSCORE, sklearnRIDGE(0.1)))
```

r_1^*	r_2^*	r_3^*	r_4^*
0.0144	0.0144	0.3583	0.3573

Table 5: Scikit-learn cross validation with $\alpha = 0.1$

As we can see in Table 5 we obtain the same results, exploiting the consistency of our work.

8 Conclusions

We used a dataset of houses in USA to build, train and test a model in order to infer housing prices. During the experiments we showed how the best value for α tends to be close to 0, this means that the amount of regularization needed to reduce the risk of overfitting is dramatically small, meaning that there is no instability.

9 Future works

Obviously, this project can be extended and deepened, for example trying other approaches like data augmentation, analysing if there are optimal values for k such that better k-fold cross validation is performed or, since no particular regularization here was needed, simply creating a linear regression model to see how different results are.

We put these points as possible guidelines for future works about this learning task.

References

- [1] M Abdalla and Khaled IA Almghari. “Remedy of multicollinearity using Ridge regression”. In: *Journal of Al Azhar University-Gaza (Natural Sciences)* 13 (2011), pp. 119–134.
- [2] Paul Allison. *When can you safely ignore multicollinearity?* URL: <https://statisticalhorizons.com/multicollinearity>. Sept. 2012.
- [3] Nicolò Cesa-Bianchi. *Statistical methods for machine learning, linear prediction*. URL: <http://cesa-bianchi.di.unimi.it/MSA/Notes/linear.pdf>. Apr. 2020.
- [4] Nicolò Cesa-Bianchi. *Statistical methods for machine learning, statistical risk*. URL: <http://cesa-bianchi.di.unimi.it/MSA/Notes/statRisk.pdf>. Mar. 2020.
- [5] Zakaria Jaadi. *A step by step explanation of principal component analysis*. URL: <https://builtin.com/data-science/step-step-explanation-principal-component-analysis>. Sept. 2019.
- [6] Barnabás Póczos. *Advanced introduction to machine learning CMU-10715, risk minimization*. URL: <https://www.cs.cmu.edu/~epxing/Class/10715/lectures/RiskMin.pdf>. Nov. 2015.
- [7] Wessel N van Wieringen. “Lecture notes on ridge regression”. In: *arXiv preprint arXiv:1509.09169* (2015).