

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Objectives</b>	<b>3</b>
2.1	Explore issues that arise from singularity. . . . .	3
2.2	Understand singular value decomposition (SVD). . . . .	3
2.3	Investigate how software deals with singularity when building linear models.	3
<b>3</b>	<b>Methodology</b>	<b>4</b>
<b>4</b>	<b>Results and discussion</b>	<b>4</b>
4.1	Python's <code>scikit-learn</code> library . . . . .	4
4.2	Python's <code>statsmodel</code> library . . . . .	4
4.3	R's <code>lm</code> function . . . . .	5
4.4	Impact of multicollinearity . . . . .	5
4.5	Methods to detect collinearity . . . . .	5
4.6	Confounding variables . . . . .	6
<b>5</b>	<b>Conclusion</b>	<b>6</b>
<b>6</b>	<b>Future recommendations</b>	<b>6</b>
<b>7</b>	<b>Appendix</b>	<b>6</b>
<b>8</b>	<b>References</b>	<b>6</b>

# 1 Introduction

Linear regression is widely used in many applications to model linear relationships between a set of predictors and a response. One of the more common and simpler methods of fitting such a model is known as the ordinary least squares method (OLS), which is found to be a Best Linear Unbiased Estimator (BLUE) by the Gauss-Markov theorem, that is, it gives the unbiased estimates of each regression coefficient with the lowest possible variance, provided that assumptions are satisfied. Furthermore, it has a main advantage of it being easily interpretable as a parametric model as opposed to non-parametric models.

The motivation for this project arises from one of these assumptions, that is, the data matrix,  $X$ , is full rank, i.e. none of the predictors are (nearly) perfectly correlated with each other. In cases where this is not true, we find that OLS fails to acquire any results. In such cases, there are several methods we can employ to address this issue, whether to drop one of the correlated variables, or to combine them, etc., based on the discretion of the data analyst. In this project, we investigate how two of the commonly used software for regression analysis, Python (`sklearn` and `statsmodel`) and R each deal with the collinearity problem.

Furthermore, we embrace this opportunity to explore further on the topic of multicollinearity, such as its impact on the accuracy of model prediction and model interpretation. We have also identified several known methods of detecting collinearity in the data, whether through hypothesis testing, or by using measures such as the variance inflation factor (VIF) and condition numbers/indices. In addition to aforementioned methods of addressing collinearity, shrinkage methods such as ridge regression and lasso were also considered to alleviate the impact of overfitting (due to collinearity) on variance at the cost of negligible bias in order to obtain an improved model for prediction. Another potential pitfall in regression analysis, that is confounding variables, were also briefly explored by examining examples of Simpson's paradox.

It should be noted that there are different conventions to what collinearity and multicollinearity mean according to varying sources. In this report, we will use these two terms interchangeably.

## 2 Objectives

There are three main objectives for this project:

### 2.1 Explore issues that arise from singularity.

In practical settings, it is unlikely to obtain a set of perfectly correlated data due to inevitable noise. However, even near perfect collinearity can lead to singularity and subsequently the failure of OLS when employing the use of software. This project aims to explore how a rank-deficient data matrix can lead to singularity, how (nearly) perfect collinearity can lead to the failure of OLS, and the pitfalls caused by moderate to high multicollinearity in general in terms of model interpretation, model prediction, and model validity. Subsequently, we look into commonly used methods to handle situations where regressors are highly correlated depending on purpose of analysis.

### 2.2 Understand singular value decomposition (SVD).

The least squares method estimates the regression coefficients of a model with  $p$  predictors by fitting it to a dataset with  $n$  observations using the following equation

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y},$$

where  $\beta$  is a  $p \times 1$  column matrix,  $\mathbf{X}$  is a  $n \times p$  data matrix, and  $\mathbf{y}$  is a  $n \times 1$  column matrix containing the observed response values. With the presence of perfect collinearity, the matrix  $\mathbf{X}^T \mathbf{X}$  is singular and hence  $(\mathbf{X}^T \mathbf{X})^{-1}$  does not exist. However, the least squares solver of Python uses the Moore-Penrose inverse (or pseudoinverse) which gives the best approximation of  $(\mathbf{X}^T \mathbf{X})^{-1}$ , and can be easily computed using the SVD of  $\mathbf{X}^T \mathbf{X}$ .

### 2.3 Investigate how software deals with singularity when building linear models.

It is observed that there exists discrepancy between how different software handles highly correlated predictors when fitting a multiple linear regression model via OLS. In this project, we aim to understand the reason behind this difference in behaviour by referring to standard documentation and examining program output.

### 3 Methodology

We narrow our focus to three least squares solvers of interests: Python's `scikit-learn` library, Python's `statsmodel` library, and R's `lm` function. In each case, we have randomly generated several sets of artificial data: one set using two correlated regressors and one set using three regressors where only one pair is correlated, each with 100 samples. The degree of collinearity is varied between high and perfect to observe differences in output. The pseudoinverse of the data matrix is also computed and multiplied with the generated response to give the theoretical minimum norm solution for the coefficient estimates and the results are compared. The use of artificial data has the advantage where the true relationship between the response and predictors is known and self-determined and that the degree of correlation and noise in the data can be manually specified.

We also performed regression analysis on built-in datasets such as California housing and diabetes dataset in `scikit-learn` and tried fitting highly correlated models. Subsequently, ridge regression models and lasso models were also used to examine its effectiveness in increasing the  $R^2$  score of these models.

## 4 Results and discussion

### 4.1 Python's `scikit-learn` library

(include screenshot of results?)

- Evenly distributes weights amongst the correlated predictors (found out to be due to `StandardScaler`)
- Does not give flags or warnings upon detection of high multicollinearity. Needs to be assessed manually by the data analyst.

### 4.2 Python's `statsmodel` library

(include screenshot of results?)

- Looks like it evenly distributes weights throughout the correlated predictors.
- Apparently uses pseudoinverse, and that the results obtained don't carry meaning.
- Gives warning upon detect of high to perfect collinearity.

### 4.3 R's `lm` function

(include screenshot of results?)

- Looks like it dropped one collinear column and add the two coefficients together.
- Seems to combine collinear columns.
- Gives a warning when highly correlated variables are detected
- Option available to abort process upon detection of high multicollinearity.

### 4.4 Impact of multicollinearity

- inflate variance and covariances of OLS estimates
- important regressors may have low t-stats
- OLS estimates and their variance become sensitive to minor changes in data
- removing/adding samples changes the regressors chosen by variable selection
- inflate some of the OLS estimates
- incorrect sign of OLS estimate

### 4.5 Methods to detect collinearity

One simple way to detect collinearity is through a test for correlation, where we test the hypothesis

$$H_0 : \rho = 0 \quad \text{vs.} \quad H_1 : \rho \neq 0,$$

where  $\rho$  is the population correlation coefficient between a chosen pair of predictors. However, there are two main reasons to opt for an alternative method:

1. as the number of predictors included in the model increase, the number of tests required to cover every combination of pairs increases. With large number of tests, the probability that we commit at least one error is high and hence we risk incorrectly identifying correlated pairs;
2. even if we find that every pair of predictors has no significant correlation, we cannot conclude the same about predictors amongst groups of other numbers.

Therefore, we turn to a more commonly used method of identifying highly correlated predictors: the variance inflation factor (VIF). (Formula and interpretation of formula and VIF values).

(Condition numbers)

## 4.6 Confounding variables

In addition to the pitfall of collinearity in performing regression analysis, confounding variables, or confounders, can also severely distort our results if not given attention. The concept of confounders is illustrated using Simpson's paradox, which states that

## 5 Conclusion

Difference in flag/no flag: indicates general attitude of different fields that uses regression modelling. ML engineers who are more concerned about making accurate predictions would not be as concerned about the impacts of collinearity on model interpretability given that the fitted model returns a satisfactory test score. This is in stark contrast to statisticians who are more interested in investigating the relationships between the predictors and the response. Hence, they are more likely to use `statsmodel` or R as opposed to `scikit_learn` which in addition to flagging high multicollinearity, also affords a more extensive function library for the purpose of performing analysis on the model.

## 6 Future recommendations

1. Explore the same issue in other software commonly used in the industry such as (?)
2. Confounding and causality.

## 7 Appendix

## 8 References

Citations to be added in second draft.