

Contents

1	Introduction	2
2	Objectives	3
3	Methodology	4
3.1	Methods to detect collinearity	5
3.1.1	Variance inflation factor (VIF)	5
3.1.2	Condition numbers	6
3.2	Methods to address collinearity	6
4	Results and discussion	7
4.1	Python's <code>scikit-learn</code> library	7
4.2	Python's <code>statsmodel</code> library	8
4.3	R's <code>lm</code> function	8
4.4	Impact of multicollinearity	8
4.5	Confounding variables	9
5	Conclusion	10
6	Future recommendations	10
A	Singular value decomposition	11
A.1	Prerequisites	11
A.2	Singular value decomposition	11
A.3	Computation of the SVD	13
B	Referenced results	13
B.1	Fitting the OLS model	13
B.2	Minimum norm approximation using pseudoinverse	13
B.3	Singularity due to perfect collinearity	14
B.4	Effect of <code>StandardScaler</code> on collinear columns in data matrix	16

1 Introduction

Multiple 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, \mathbf{X} , is full rank, i.e. none of the predictors are (nearly) perfectly correlated with each other. In cases where this is not true, OLS is unable to estimate regression parameters due to singularity. 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 (`scikit_learn` and `statsmodel`) and R each handles perfect collinearity when fitting the model.

Furthermore, we embrace this opportunity to explore further on the topic of multicollinearity, such as its impact on model prediction and interpretation. We also apply 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, we consider shrinkage methods such as ridge regression and lasso to alleviate the impact of overfitting (due to collinearity) on variance at the cost of negligible bias. We also briefly explore another potential pitfall in regression analysis, that is confounding variables, by examining examples of Simpson's paradox.

Here, we lay down the general setup for fitting a linear regression model. Given an $n \times p$ data matrix \mathbf{X} (n samples and p predictors), corresponding to an $n \times 1$ column matrix \mathbf{y} containing the observed response, we model the relationship between p predictors, X_i and the response variable, Y as

$$Y = \beta_0 + \sum_{i=1}^p \beta_i X_i + \varepsilon,$$

where ε represents the random error. We fit the model to the data collected for the expected response $E[Y]$ by estimating each β_i using the least squares method to obtain $\hat{\beta}_i$ so that we

have

$$E[Y] = \hat{\beta}_0 + \sum_{i=1}^p \hat{\beta}_i X_i.$$

Finally, 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 to refer to the correlation between two or more predictors in a regression model.

2 Objectives

There are three main objectives for this project:

1. Explore issues that arise from singularity.
2. Understand singular value decomposition (SVD).
3. Investigate how software deals with singularity when building linear models.

In practical settings, it is unlikely to obtain a set of perfectly correlated data due to inevitable noise. However, even near perfect collinearity leads to singularity when attempting to fit the regression model using software. We hence aim to explore how a rank-deficient data matrix can lead to singularity and the issues that arise from it. We also explore pitfalls caused by moderate to high multicollinearity in terms of interpretation and prediction. Subsequently, we look into commonly used methods to handle situations where regressors are highly correlated depending on purpose of analysis.

The next objective is to understand how to apply the singular value decomposition (SVD) in the computations of the least square estimates. 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

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

where $\hat{\beta}$ is a $p \times 1$ column matrix containing the least square estimates, \mathbf{X} is a $n \times p$ data matrix, and \mathbf{y} is a $n \times 1$ column matrix containing the observed response values. The derivation of this formula is included in the appendix (Section B.1). 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) of \mathbf{X} , \mathbf{X}^+ , which can be used to calculate the minimum norm approximation of $(\mathbf{X}^T \mathbf{X})^{-1}$, and can be computed using the SVD of $\mathbf{X}^T \mathbf{X}$ (see appendix Section B.2). A brief section in the

appendix (Section A) is dedicated to the formulation and computation of the SVD.

Lastly, 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 in Python and R. As an extension, methods for detecting and addressing multicollinearity using software are also explored.

3 Methodology

We narrowed 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 randomly generated several sets of artificial data using three regressors where only one pair is correlated, each with 100 samples. 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. Note that the distribution used in data generation is irrelevant for our topic of study.

The data ranges from 0 to 100 and is uniformly distributed. The degree of collinearity was then varied between perfect and moderate to observe differences in output. The pseudoinverse of the data matrix was also computed and multiplied with the generated response to give the theoretical minimum norm solution for the coefficient estimates. Then, these results were compared with those given by the least squares solver function to confirm whether these functions use the pseudoinverse of the data matrix to estimate the least squares coefficients as suggested by documentation.

The true underlying relationship to which the artificial data was generated was set to be

$$Y = X_1 + 2X_2 + 3X_3,$$

where X_1 and X_2 are correlated with $X_2 = 2X_1$. We then vary the amount of uniformly distributed noise added to different generated sets: from no noise (perfect collinearity) to some number sampled from intervals $[-25,25]$, $[-50,50]$, and $[-100,100]$. The train-test ratio chosen for splitting the dataset into training and testing sets is fixed at 80-20. Lastly, to evaluate prediction performance, metrics such as the (adjusted) multiple coefficient of determination R^2 and mean squared error were used.

The scaling effect of `scikit-learn`’s `StandardScaler` on the ridge regression and

LASSO models in the presence of collinearity was also given attention. This was done by referring to standard documentation as well as by testing it using artificially generated data. For data with and without high multicollinearity present respectively, two models were each fitted: one with predictors scaled by `StandardScaler` and one without. The output was then compared.

We have also performed regression analysis on California housing and applied methods to address collinearity using `statsmodel`. Subsequently, ridge regression models and lasso models were also used to examine its effectiveness in decreasing the standard errors of these models.

3.1 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 ρ is the population correlation coefficient between a chosen pair of predictors. However, there are two main reasons to opt for an alternative method.

Firstly, as the number of predictors included in the model increases, 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. More importantly, even if we find that every pair of predictors has no significant correlation, we cannot make the same conclusion about predictors amongst groups of other numbers.

3.1.1 Variance inflation factor (VIF)

A commonly used method for identifying highly correlated predictors is by using the variance inflation factor (VIF), given by

$$\text{VIF}_k = \frac{1}{1 - R_k^2},$$

where R_k^2 is the multiple coefficient of determination of the k^{th} predictor regressed against the other predictors, assuming that it has zero mean. Intuitively, a higher value of R_k^2 indicates a good fit, hence a strong correlation between the k^{th} regressor and some of the other regressors. Note that uncorrelated regressors would simply have a coefficient close to zero.

The motivation behind the formulation of the VIF is to calculate the ratio between the variance of a coefficient estimate if it was uncorrelated with the other regressors (baseline

case), and the actual variance. For example, a coefficient estimate having a VIF of 4 corresponds to its variance being 4 times of what it would have been if it was uncorrelated with every other predictor. In other words, VIF measures the factor by which the variance of the regression coefficients is inflated above what it would have been if $R_k^2 = 0$, hence its name.

There are different rule of thumbs for the thresholds beyond which one may consider the VIF value of a particular predictor to be high, ranging roughly from 4 to 10. There is no definitively correct answer, and we have no choice but to make decisions based on our interpretation of the VIF value. Once predictors with high VIF's are identified, we can then examine the correlation matrix to identify which the predictors which each high VIF predictor is correlated to.

3.1.2 Condition numbers

Another metric for measuring the degree of collinearity is the condition number (condition indices) of the data matrix \mathbf{X} , which indicates how sensitive the response is to changes in \mathbf{X} . Aside from using it to determine whether a model is highly collinear, the condition number of $\mathbf{X}^T \mathbf{X}$ is used to check for a potential source of large numerical error due to collinearity when computing $(\mathbf{X}^T \mathbf{X})^{-1}$ using software. In this case, it indicates the multiplicative factor to which round-off errors are amplified due to computer arithmetic being performed to finite precision.

The formula for the condition number of a $m \times n$ matrix \mathbf{A} with rank $r \leq n$, denoted $\kappa(\mathbf{A})$, is given by

$$\kappa(\mathbf{A}) = \frac{\sigma_1}{\sigma_r},$$

where σ_i are the singular values of \mathbf{A} arranged in descending order. Note that definitions will differ for varying contexts, giving rise to different condition numbers. The point of discussion here is the condition number for inversion.

Similar to the case of VIF, there is no exact threshold for a condition number to be considered high, but a rule of thumb is any condition numbers exceeding 20 indicates the matrix to be ill-conditioned.

3.2 Methods to address collinearity

We present a few commonly-used methods to address moderate to perfect collinearity.

In the case of perfect collinearity, one of the correlated predictors is completely described by the other, hence our only choice is to remove one of them. This is however not

recommended for non-perfect collinearity as we are at risk of committing omitted variable bias. In such cases, we may want to either combine some of the variables or use shrinkage methods such as ridge regression or LASSO.

There are many possible ways of combining correlated variables, but in most cases the simple approach of adding the columns up (such as what R does) gives sufficiently improved results. Note that it is important that we consider the context when doing this, that is, whether combining the two variables make sense in context. An example is a model predicting the salary of employees based on predictors such as years of experience, gender, test scores for course A, and test scores for course B. Obviously, any employee who does well on one test is expected to do well in the other, and so we may consider adding the two test scores together to form a single column "total test scores for course A and B", then the results obtained may be interpreted accordingly.

For ridge regression and LASSO, otherwise known as L_2 and L_1 regularisation respectively, we find that it successfully decreases coefficient standard errors at the cost of increasing the bias (and hence decreasing R^2 score compared to the OLS model) by a negligible amount. However, in the face of perfect collinearity, ridge and lasso does not yield interpretable coefficients, where the correlated predictors have equally distributed coefficients due to the effect of standardising the predictors beforehand (see Appendix Section B.2), as the ridge and LASSO coefficient estimates are not scale equivariant.

4 Results and discussion

4.1 Python's `scikit-learn` library

After referring to documentation, we find that the least squares solver, `LinearRegression().fit()` computes the coefficient estimates using the pseudoinverse of the data matrix, which we have verified. Under any degree of collinearity, we observed that model prediction is affected minimally, attributed by R^2 scores close to 1. However, when predictors are scaled, the ridge regression or LASSO model fitted to data with perfect collinearity results in equal distributed coefficients between correlated predictors. After further investigation, we discovered that this was caused by the behaviour of `StandardScaler` from the same library which was used to scale predictors prior to fitting the model, which we show in Section B.4.

We also note that the solver does not give any warnings of high collinearity when displaying results and does not apply any methods to handle highly collinear data.

4.2 Python's statsmodel library

(include screenshot of results?)

- Looks like it evenly distributes weights throughout the correlated predictors.
- Uses pseudoinverse
- 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

High multicollinearity mainly poses a major problem to model interpretation. Notably, it inflates the variance and covariances of least squares estimates and causes estimates and their variance to become sensitive to minor changes in data. As a result, we may obtain different results when performing variable selection after removing or adding samples from the dataset, or obtain incorrect signs of coefficients estimates. Another curious effect is that the t -values of correlated predictors may be insignificant yet removing these predictors from the model results in a lower R^2 value. Although this seems contradictory by intuition, it is actually reasonable behaviour, as the prediction is contributed by both variables while the t -values indicate the contribution of one variable after accounting for the others.

When two predictors have perfect correlation, we have what is known as perfect collinearity. In such a case, at least one of the columns of the data matrix \mathbf{X} is a linear combination of the other columns. This causes the matrix $\mathbf{X}^T \mathbf{X}$ to be singular (see appendix Section B.3), consequently causing $\hat{\beta}$ unable to be computed (see Equation 1). However, statsmodel's and scikit_learn's least square solvers compute the pseudoinverse of $\mathbf{X}^T \mathbf{X}$ when estimating the regression coefficients, giving the least norm approximated solution (which is interpretable) when \mathbf{X} is rank-deficient. However, the aforementioned impacts of high collinearity still apply hence it should be addressed.

It is important to note that collinearity does not, in any way, impact predictions or its precision. Hence, if the model is not expected to be interpreted, multicollinearity is no issue of concern and may be ignored. Moreover, even for more statistically-oriented analysts, there may also not be a need to address high collinearity provided that it does not inflate standard errors to a problematic extent (see "A caution regarding rules of thumb for variance inflation factors" by O'Brien), as there are many other factors which can suppress the inflation effect of standard errors, one of which is the sample size. Quoting the words of O'Brien, "collinearity does not hurt as long as it does not bite".

4.5 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, where the relationship between individual predictors and the response disappears or reverses when all predictors are included in a single model. The variables which cause such occurrences are referred to as confounders.

We give the classic textbook example of the relationship between ice cream sales and shark attacks. We initially observe a positive correlation between these two variables, yet upon the inclusion of the confounder, temperature, we find that this correlation vanishes. One should avoid the confusion between confounders and correlated variables by noting that in addition to being correlated, confounders should also have a causal relationship to the variables involved in the confounding effect.

Although including confounders inherently introduces collinearity into the model, failing to account for them causes us to commit omitted variable bias in addition to causing the aforementioned confounding effects according to Simpson's paradox. (We note the difference between these confounding effects and the effect of collinearity on the coefficient signs: the latter may still occur even after including confounders in the model due to high collinearity causing coefficient estimates to become sensitive to changes in observed data). Hence, by weighing each consequence, we generally always include confounders as we have well-established methods for effectively addressing collinearity. Unfortunately, it is more often the case that confounders are excluded unintentionally, as there is no standard method of identifying them. Furthermore, this is not a problem exclusive to regression modelling, but observational studies in general. Hence, we highlight this further in Section 6.

5 Conclusion

After examining Python's and R's program output, we have observed differences in how they each handle the issue of multicollinearity. In particular, R's behaviour of combining correlated predictors while giving a warning upon detection of high multicollinearity aligns best with standard practice. It is also found that in the process of estimating the regression coefficients, R's `lm`, Python's `statsmodel`, and `scikit_learn` calculates the pseudoinverse of $X^T X$ rather than $(X^T X)^{-1}$ so that a result may be obtained even in the presence of singularity.

Additionally, we have observed a difference between the three's behaviour of flagging problematic multicollinearity, which indicates the general attitude of different fields that uses regression modelling. For example, `scikit_learn` aligns with machine learning engineers who are generally more concerned about making accurate predictions and hence would not be as concerned about the impacts of collinearity on model interpretability. This is in stark contrast to statisticians who bear greater interest in investigating the relationships between the regressors and the response. Hence, as opposed to `scikit_learn`, they would be more likely to use `statsmodel` or R, which in addition to flagging high multicollinearity, also affords a more extensive function library for the purpose of performing analysis on the model.

Taking these distinctions into account, it depends ultimately on the data analyst to make the final decision on which approach is best to handle anomalies in the model after performing analysis according to context and domain knowledge. Nevertheless, it is still crucial that one is made aware of these differences so as to select the appropriate tools depending on the individual's purpose and objective of regression modelling.

6 Future recommendations

1. Explore the same issue in other software commonly used in the industry such as (?)
2. Confounding and causality.
3. Partial effect, total effect, direct effect

A Singular value decomposition

The SVD is a decomposition of any matrix into a sum of rank-1 matrices. Before we give its formulation, we first state several definitions as well as results without proof.

A.1 Prerequisites

Definition A.1 (Symmetric matrices). *An $n \times n$ square matrix \mathbf{X} is symmetric if $\mathbf{A} = \mathbf{A}^T$.*

Definition A.2 (Orthogonal matrices). *An $n \times n$ square matrix $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{pmatrix}$ is orthogonal if its columns \mathbf{x}_i are orthonormal, that is, they are orthogonal to each other and each have a norm of 1, i.e. $\mathbf{x}_i \mathbf{x}_j^T = 0$ for $i \neq j$ and $\mathbf{x}_i \mathbf{x}_i^T = 1$ for $i = j$, where $1 \leq i, j \leq n$.*

Proposition A.1 (Spectral/eigen-decomposition). *Any $n \times n$ symmetric matrix \mathbf{A} can be written as*

$$\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T = \sum_{i=1}^n \lambda_i \mathbf{q}_i \mathbf{q}_i^T,$$

where $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ is an $n \times n$ diagonal matrix with diagonals equal to the eigenvalues of \mathbf{A} and \mathbf{Q} is an orthogonal matrix whose columns are unit eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$ of \mathbf{A} .

Definition A.3 (Singular vectors and singular values). *Let \mathbf{A} be a $n \times p$ matrix, then σ is a singular value with corresponding left and right singular vectors \mathbf{u} and \mathbf{v} respectively if*

$$\mathbf{A} \mathbf{v} = \sigma \mathbf{u} \quad \text{and} \quad \mathbf{A}^T \mathbf{u} = \sigma \mathbf{v}.$$

A.2 Singular value decomposition

Proposition A.2 (Singular value decomposition). *Let \mathbf{A} be an $n \times p$ matrix with rank r , where $1 \leq r \leq \min(n, p)$. Then there exists an $n \times r$ matrix $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_r\}$, an $p \times r$ matrix $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_r\}$, and an $r \times r$ matrix $\mathbf{\Sigma} = \text{diag}\{\sigma_1, \dots, \sigma_r\}$, such that*

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T,$$

where $\mathbf{U} \mathbf{U}^T = \mathbf{I}_r = \mathbf{V} \mathbf{V}^T$, and $\sigma_1 \geq \dots \geq \sigma_r > 0$.

The form given above is called the compact SVD. The non-compact form is given by

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T,$$

where \mathbf{U} is an $n \times n$ orthogonal matrix, \mathbf{V} is a $p \times p$ orthogonal matrix and $\mathbf{\Sigma}$ is a $n \times p$ diagonal matrix, leaving the remaining diagonals after the r^{th} entry zero.

Proposition A.3. *Let A be an $n \times p$ matrix with rank r , then*

$$\text{rank}(A^T A) = \text{rank}(A) = r.$$

Proof. To prove this, we show that A and $A^T A$ have the same null space and hence the same nullity. We do this by proving that $\text{nullity}(A) \subseteq \text{nullity}(A^T A)$ and $\text{nullity}(A^T A) \subseteq \text{nullity}(A)$ separately.

For the first case, let $x \in \mathbb{R}^p$ such that $Ax = \mathbf{0}_n$. Then, $A^T Ax = A^T(\mathbf{0}_n) = \mathbf{0}_r$. So $\text{nullity}(A) \subseteq \text{nullity}(A^T A)$. For the second case, $x \in \mathbb{R}^p$ such that $A^T Ax = \mathbf{0}$. Left multiplying both sides by x^T , we have

$$x^T A^T Ax = (Ax)^T (Ax) = \|Ax\|^2 = 0,$$

which implies that $Ax = \mathbf{0}$, so $\text{nullity}(A^T A) \subseteq \text{nullity}(A)$.

Then since the dimension of the domain of $A^T A$ and A are both p , by the rank-nullity formula, they have the same rank. \square

Proposition A.4. *Let A be an $n \times p$ matrix of rank r . Then the non-zero eigenvalues of both AA^T and $A^T A$ are $\sigma_1^2, \dots, \sigma_r^2$. Furthermore, the corresponding unit eigenvectors of AA^T and $A^T A$ are given by the columns of U and V respectively.*

Proof. We prove this result by first noting that $A^T A$ is a $p \times p$ symmetric matrix, so we can write its spectral decomposition as

$$A^T A = V \Lambda V^T$$

where V is a $p \times r$ semi-orthogonal matrix containing the orthonormal eigenvectors of $A^T A$ (and hence its columns v_i are orthonormal) and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_r)$ is a diagonal matrix with its diagonals containing eigenvalues $\lambda_1 \geq \dots \geq \lambda_r > 0$. Note that we know that there are r eigenvalues by Proposition A.3.

Now, we let each $\sigma_i = \sqrt{\lambda_i}$ and $u_i = \frac{1}{\sigma_i} A v_i$ for $i = 1, \dots, r$. We then show that the vectors u_i are orthonormal and subsequently, u_i and v_i are left and right singular vectors corresponding to singular values σ_i as follows:

$$u_i^T u_j = \frac{1}{\sigma_i \sigma_j} v_i^T A^T A v_j = \frac{1}{\sigma_i \sigma_j} v_i^T (\lambda_j v_j) = \frac{\sigma_j^2}{\sigma_i \sigma_j} v_i^T v_j = \frac{\sigma_j^2}{\sigma_i \sigma_j} v_i \cdot v_j.$$

Notice that if $i = j$, then the expression simplifies to just $\|v_i\|^2 = 1$ since the vectors v_i are orthonormal. If $i \neq j$, then the dot product equals to zero, implying that the vectors u_i

are orthonormal. Now, following the definition of singular vectors, we consider

$$A^T \mathbf{u}_i = \frac{1}{\sigma_i} (A^T A \mathbf{v}_i) = \frac{\sigma_i^2}{\sigma_i} \mathbf{v}_i = \sigma_i \mathbf{v}_i,$$

so \mathbf{v}_i are right singular vectors corresponding to singular values σ_i while \mathbf{u}_i by our formulation are left singular vectors. Now all that is left is to construct $\mathbf{U} = \begin{pmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_r & \dots & \mathbf{u}_n \end{pmatrix}$ and $\mathbf{\Sigma} = \text{diag}\{\sigma_1, \dots, \sigma_r, 0, \dots, 0\}$. \square

A.3 Computation of the SVD

Let A be any $m \times n$ matrix. To compute the SVD of A , we first solve for eigenvalues λ of AA^T or $A^T A$ depending on which results in a matrix of smaller size. Then we take the singular values $\sigma = \sqrt{\lambda}$ and construct $\mathbf{\Sigma}$. Since in regression modelling, we often have more samples than predictors (which translates to more rows than columns), we choose $A^T A$. Now, we solve for right singular vectors by using its definition:

$$(A^T A - \lambda I) \mathbf{v} = \mathbf{0}.$$

Then, we convert the right singular vectors to unit vectors (so that \mathbf{U} is orthogonal) and construct matrix \mathbf{U} . Compute \mathbf{V} by using the definition of right singular vectors:

$$\sigma_i \mathbf{v}_i = A^T \mathbf{u}_i.$$

B Referenced results

B.1 Fitting the OLS model

(explain using the concept of orthogonal projection)

B.2 Minimum norm approximation using pseudoinverse

We can compute the pseudoinverse of $X^T X$ using its SVD as follows.

Let $A = X^T X$. If the SVD of $A = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, then the pseudoinverse of A , A^+ , is given

by

$$\begin{aligned}
A^+ &= ((\mathbf{U}\Sigma\mathbf{V}^T)^T)(\mathbf{U}\Sigma\mathbf{V}^T)^{-1}(\mathbf{U}\Sigma\mathbf{V}^T)^T \\
&= (\mathbf{V}\Sigma\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T)^{-1}(\mathbf{U}\Sigma\mathbf{V}^T)^T \\
&= (\mathbf{V}\Sigma^2\mathbf{V}^T)^{-1}(\mathbf{V}\Sigma\mathbf{U}^T) \\
&= (\mathbf{V}^T)^{-1}(\Sigma^{-1})^2(\mathbf{V})^{-1}\mathbf{V}\Sigma\mathbf{U}^T \\
&= \mathbf{V}\Sigma^{-1}\mathbf{U}^T
\end{aligned}$$

Now, for $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, if the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ has solutions, then $\mathbf{x}^* = \mathbf{A}^+\mathbf{b}$ is an exact solution and has the smallest possible norm, i.e. $\|\mathbf{x}^*\| \leq \|\mathbf{x}\|$ for all \mathbf{x} .

The proof is as follows: Since \mathbf{A}^+ is a generalised inverse, $\mathbf{A}^+\mathbf{b}$ must be a solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$. Now, for any solution $\mathbf{x} \in \mathbb{R}^n$, consider its orthogonal decomposition via $\mathbf{A}^+\mathbf{A} \in \mathbb{R}^{n \times n}$:

$$\mathbf{x} = (\mathbf{A}^+\mathbf{A})\mathbf{x} + (\mathbf{I} - \mathbf{A}^+\mathbf{A})\mathbf{x} = \mathbf{A}^+\mathbf{b} + (\mathbf{I} - \mathbf{A}^+\mathbf{A})\mathbf{x}.$$

Then by the Pythagorean theorem, we have

$$\|\mathbf{x}\|^2 = \|\mathbf{A}^+\mathbf{b}\|^2 + \|(\mathbf{I} - \mathbf{A}^+\mathbf{A})\mathbf{x}\|^2 \geq \|\mathbf{A}^+\mathbf{b}\|^2.$$

Hence $\|\mathbf{x}\| \geq \|\mathbf{A}^+\mathbf{b}\|$.

B.3 Singularity due to perfect collinearity

Let $\mathbf{X} = \begin{pmatrix} \mathbf{u} & k\mathbf{u} & \mathbf{v} \end{pmatrix}$ where $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ are arbitrarily linearly independent column vectors containing n observations of a particular feature with $\mathbf{u}^T = \begin{pmatrix} u_1 & \dots & u_n \end{pmatrix}$ and $\mathbf{v}^T = \begin{pmatrix} v_1 & \dots & v_n \end{pmatrix}$, with $k \in \mathbb{R}$ such that the column $k\mathbf{u}$ is a scalar multiple of the first column.

Consider

$$\begin{aligned}
X^T X &= \begin{pmatrix} \mathbf{u}^T \\ k\mathbf{u}^T \\ \mathbf{v}^T \end{pmatrix} \begin{pmatrix} \mathbf{u} & k\mathbf{u} & \mathbf{v} \end{pmatrix} \\
&= \begin{pmatrix} u_1 & \dots & u_n \\ ku_1 & \dots & ku_n \\ v_1 & \dots & v_n \end{pmatrix} \begin{pmatrix} u_1 & ku_1 & v_1 \\ \vdots & \vdots & \vdots \\ u_n & ku_n & v_n \end{pmatrix} \\
&= \sum_{i=1}^n \begin{pmatrix} u_i \\ ku_i \\ v_i \end{pmatrix} \begin{pmatrix} u_i & ku_i & v_i \end{pmatrix} \\
&= \sum_{i=1}^n \begin{pmatrix} u_i^2 & ku_i^2 & u_i v_i \\ ku_i^2 & k^2 u_i^2 & ku_i v_i \\ u_i v_i & ku_i v_i & v_i^2 \end{pmatrix}.
\end{aligned}$$

The second row of $X^T X$ is k times the first row for every i . After summing up the rank-1 matrices, the second row of $X^T X$ will also be k times the first. So, $X^T X$ is not full rank, and hence is singular.

By writing matrix-matrix products in this form, each term in the sum is a rank-1 matrix. However, under no collinearity, the rows in every term will not necessarily be dependent on some other rows **by the same factor**. In the formulation above, the third row differs from the first row by a multiplicative factor of v_i/u_i (this is consistent with the case of collinear columns since $ku_i/u_i = k$), which is a constant under no collinearity. This fraction will differ for varying i , thus after adding each term up, there is no common factor between the two rows.

Alternatively for a simpler and less analytic approach, we look at each entry in the product matrix as multiplying rows with columns:

$$X^T X = \begin{pmatrix} \mathbf{u} \cdot \mathbf{u} & \mathbf{u} \cdot (k\mathbf{u}) & \mathbf{u} \cdot \mathbf{v} \\ (k\mathbf{u}) \cdot \mathbf{u} & (k\mathbf{u}) \cdot (k\mathbf{u}) & (k\mathbf{u}) \cdot \mathbf{v} \\ \mathbf{v} \cdot \mathbf{u} & \mathbf{v} \cdot (k\mathbf{u}) & \mathbf{v} \cdot \mathbf{v} \end{pmatrix} = \begin{pmatrix} \mathbf{u} \cdot \mathbf{u} & k\mathbf{u} \cdot \mathbf{u} & \mathbf{u} \cdot \mathbf{v} \\ k(\mathbf{u} \cdot \mathbf{u}) & k(k\mathbf{u} \cdot \mathbf{u}) & k(\mathbf{u} \cdot \mathbf{v}) \\ \mathbf{v} \cdot \mathbf{u} & k(\mathbf{v} \cdot \mathbf{u}) & \mathbf{v} \cdot \mathbf{v} \end{pmatrix}.$$

We observe that the second row is k times the first row as expected. The third row again is not collinear with the first two. Suppose they are, then we solve for some $\lambda \in \mathbb{R}$ such that

$$\mathbf{v} \cdot \mathbf{u} = \lambda(\mathbf{u} \cdot \mathbf{u}) \quad \text{and} \quad \mathbf{v} \cdot \mathbf{v} = \lambda(\mathbf{u} \cdot \mathbf{v}).$$

It then follows that

$$\begin{aligned}
\|v\|^2 &= v \cdot v \\
&= \lambda(v \cdot u) \\
&= \lambda^2(u \cdot u) \\
&= \lambda^2\|u\|^2 \\
\|v\| &= \lambda\|u\| \\
&= \frac{v \cdot u}{\|u\|^2}\|u\| \\
v &= \frac{v \cdot u}{\|u\|^2}u
\end{aligned}$$

i.e. v is its own orthogonal projection on vector u , implying that they are collinear, which is a contradiction. Hence the third row is linearly independent of the first two rows.

Since the least squares coefficient estimates are given by

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

the singularity of $X^T X$ causes least squares to fail to obtain estimates for .

B.4 Effect of StandardScaler on collinear columns in data matrix

From the documentation, `scikit_learn`'s `StandardScaler` calculates the Z-score of each entry by subtracting each entry by the mean of its column and dividing by the standard deviation of its column. Suppose that the true relationship between regressors x_1 and x_2 is $x_2 = kx_1$. Then we have

$$\bar{x}_2 = k\bar{x}_1 \quad \text{and} \quad s_2^2 = k^2 s_1^2.$$

So for collinear columns x_1 and x_2 , the entries of scaled columns u_1 and u_2 are

$$u_{i1} = \frac{x_{i1} - \bar{x}_1}{s_1} \quad \text{and} \quad u_{i2} = \frac{x_{i2} - \bar{x}_2}{s_2} = \frac{kx_{i1} - k\bar{x}_1}{ks_1} = \frac{x_{i1} - \bar{x}_1}{s_1} = u_{i1}.$$

References