

Pizza Example: Pizza sales are related to both the number of ads and also the cost x_1 of the ads. This became evident from the two SLRs:

$$y = \beta_0 + \beta_1 x_1 + \epsilon \quad y = \beta_0 + \beta_2 x_2 + \epsilon$$

$H_0: \beta_1 = 0$ was rejected $H_0: \beta_2 = 0$ was rejected

But x_1 and x_2 were also strongly linearly related. A consequence of this was that neither x_1 nor x_2 seemed important in the context of a model that contained both of them.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$H_0: \beta_1 = 0 \quad H_0: \beta_2 = 0$$

these were not rejected

Thus, the full model is misleading us into think x_1 and x_2 are unimportant when they really are. This is happening because of variance inflation.

Credit Card Example: $y = \text{credit card balance}$

$x_1 = \text{credit card holder's age}$

$x_2 = \text{credit card holder's credit limit}$

$x_3 = \text{credit card holder's credit rating}$

Like in the Pizza example, a strong linear relationship between x_2 and x_3 confused their importance.

Detecting Multicollinearity

Multicollinearity between two explanatory variables can easily be identified with

- Correlation Matrix
- Scatterplot Matrix

Since these will clearly identify strong pairwise linear relationships among the explanatory variables.

However this approach will not (as easily) identify multicollinearity between more than two explanatory variables. To identify this type of multicollinearity, we'll use variance inflation factors (VIF):

$$\text{VIF}_j = \frac{\text{Var}[\beta_j]}{\text{Var}[\beta_j^*]} \geq 1$$

where β_j^* is the coefficient of a simple linear regression including only x_j , and β_j is the coefficient of x_j in a larger model containing other explanatory variables.

Thus, if VIF_j is large then we have evidence that x_j is part of a multicollinearity problem.

Another way to diagnose this problem would be to fit a linear regression where x_j is treated as the response and all of the remaining x 's are the explanatory variables:

$$x_j = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_{j-1} x_{j-1} + \alpha_{j+1} x_{j+1} + \dots + \alpha_p x_p + \epsilon$$

Such a model is fit for every $j=1, 2, \dots, p$ and each time R_j^2 is calculated (i.e., the coefficient of determination as calculated from the model with x_j as the response). If R_j^2 is large, it suggests that x_j is strongly related to (some of) the other explanatory variables.

It can be shown that

$$\text{Var}[\beta_j^*] \left(\frac{1}{1-R_j^2} \right) = \text{Var}[\beta_j]$$

Values of VIF_j much larger than 1 provide evidence of multicollinearity. In particular $\text{VIF}_j \geq 5$ or $\text{VIF}_j \geq 10$ indicate serious multicollinearity

$$R_j^2 \geq 0.8$$

$$R_j^2 \geq 0.9$$

To eliminate multicollinearity, calculate the VIF for every explanatory variable and focus attention on the ones that are greater than 5 or 10. Among these ones, eliminate the explanatory variable with the largest VIF (since this means it is most strongly predicted by the others). We eliminate this one because it is redundant.

This procedure should be done repetitively until no more multicollinearity exists.

Pizza Example: $\text{VIF}_1 = \text{VIF}_2 = 5.339243 > 5$

$\therefore x_1$ and x_2 are multicollinear and one should be removed from the model. Since $\text{VIF}_1 = \text{VIF}_2$ it doesn't matter which one is removed.

Credit Card Example: $\text{VIF}_1 = 1.011385 \quad \text{VIF}_2 = 160.592880 \quad \text{VIF}_3 = 160.648301$

$\therefore x_2$ and x_3 are multicollinear and one should be removed from the model. Since $\text{VIF}_3 > \text{VIF}_2$, x_3 should be removed.

Model Selection

The idea with model selection is that given q potential explanatory variables v_1, v_2, \dots, v_q , we want to find the best subset x_1, x_2, \dots, x_p ($p \leq q$) that provides the best model

• best in terms of "fit"

• best in terms of predictive performance.

Thus we want to select $(x_1, x_2, \dots, x_p) \subseteq (v_1, v_2, \dots, v_q)$.

such that

- no important variable is left out of the model
- no unimportant variable is included in the model.

Ideally we'd like an automated way of doing this.