The modified version is known as the **centered** $R^2$, and we will denote it by $R_c^2$. It is defined as

$$R_c^2 \equiv 1 - \frac{\|\boldsymbol{M}_X\boldsymbol{y}\|^2}{\|\boldsymbol{M}_\iota\boldsymbol{y}\|^2}, \tag{1.09}$$

where

$$\boldsymbol{M}_\iota \equiv \mathbf{I} - \boldsymbol{\iota}(\boldsymbol{\iota}^\top\boldsymbol{\iota})^{-1}\boldsymbol{\iota}^\top = \mathbf{I} - n^{-1}\boldsymbol{\iota}\boldsymbol{\iota}^\top$$

is the matrix that projects off the space spanned by the constant vector $\boldsymbol{\iota}$, which is simply a vector of $n$ ones. When any vector is multiplied by $\boldsymbol{M}_\iota$, the result is a vector of deviations from the mean. Thus what the centered $R^2$ measures is the proportion of the total sum of squares of the regressand *around its mean* that is explained by the regressors.

An alternative expression for $R_c^2$ is

$$\frac{\|\boldsymbol{P}_X\boldsymbol{M}_\iota\boldsymbol{y}\|^2}{\|\boldsymbol{M}_\iota\boldsymbol{y}\|^2}, \tag{1.10}$$

but this is equal to (1.09) only if $\boldsymbol{P}_X\boldsymbol{\iota} = \boldsymbol{\iota}$, which means that $\mathcal{S}(\boldsymbol{X})$ must include the vector $\boldsymbol{\iota}$ (so that either one column of $\boldsymbol{X}$ must be a constant, or some linear combination of the columns of $\boldsymbol{X}$ must equal a constant). In this case, the equality must hold, because

$$\boldsymbol{M}_X\boldsymbol{M}_\iota\boldsymbol{y} = \boldsymbol{M}_X(\mathbf{I} - \boldsymbol{P}_\iota)\boldsymbol{y} = \boldsymbol{M}_X\boldsymbol{y},$$

the second equality here being a consequence of the fact that $\boldsymbol{M}_X$ annihilates $\boldsymbol{P}_\iota$ when $\boldsymbol{\iota}$ belongs to $\mathcal{S}(\boldsymbol{X})$. When this is not the case and (1.10) is not valid, there is no guarantee that $R_c^2$ will be positive. After all, there will be many cases in which a regressand $\boldsymbol{y}$ is better explained by a constant term than by some set of regressors that does not include a constant term. Clearly, if (1.10) is valid, $R_c^2$ must lie between 0 and 1, since (1.10) is then simply the uncentered $R^2$ for a regression of $\boldsymbol{M}_\iota\boldsymbol{y}$ on $\boldsymbol{X}$.

The use of the centered $R^2$ when $\boldsymbol{X}$ does not include a constant term or the equivalent is thus fraught with difficulties. Some programs for statistics and econometrics refuse to print an $R^2$ at all in this circumstance; others print $R_u^2$ (without always warning the user that they are doing so); some print $R_c^2$, defined as (1.09), which may be either positive or negative; and some print still other quantities, which would be equal to $R_c^2$ if $\boldsymbol{X}$ included a constant term but are not when it does not. Users of statistical software, be warned!

Notice that $R^2$ is an interesting number only because we used the least squares estimator $\hat{\boldsymbol{\beta}}$ to estimate $\boldsymbol{\beta}$. If we chose an estimate of $\boldsymbol{\beta}$, say $\tilde{\boldsymbol{\beta}}$, in any other way, so that the triangle in Figure 1.3 were no longer a right-angled triangle, we would find that the equivalents of the two definitions of $R^2$, (1.09) and (1.10), were not the same:

$$1 - \frac{\|\boldsymbol{y} - \boldsymbol{X}\tilde{\boldsymbol{\beta}}\|^2}{\|\boldsymbol{y}\|^2} \neq \frac{\|\boldsymbol{X}\tilde{\boldsymbol{\beta}}\|^2}{\|\boldsymbol{y}\|^2}.$$