$$SST = SSR + SSE$$

↑

fixed! only a function of $y$ and not $X$!

$$\Rightarrow SSR \uparrow \Rightarrow SSE \downarrow \Rightarrow R^2 \uparrow, \; RMSE \downarrow$$
$$SSR \downarrow \Rightarrow SSE \uparrow \Rightarrow R^2 \downarrow, \; RMSE \uparrow$$

$$\hat{\vec{y}} = H\vec{y} = QQ^T\vec{y} = \sum_{j=0}^{p} proj_{\vec{e}_j}(\vec{y})$$

since each $\vec{e}_j$ is orthogonal $\Rightarrow proj_{\vec{e}_j}(\vec{y})$ are all orthogonal since $\exists c_j \; \forall \vec{y}\ldots \; proj_{\vec{e}_j}(\vec{y}) = c_j\vec{e}_j$

$\Rightarrow$ by Pythag. Thm. in $p+1$ dimensions...

$$\|\hat{\vec{y}}\|^2 = \sum_{j=0}^{p} \|proj_{\vec{e}_j}(\vec{y})\|^2$$

$$SSR = \sum_{i=1}^{n}(\hat{y}-\bar{y})^2 = \sum_{i=1}^{n}\hat{y}_i^2 - \sum_{i=1}^{n}2\hat{y}_i\bar{y} + \sum_{i=1}^{n}\bar{y}^2$$

$$= \|\hat{\vec{y}}\|^2 - 2\bar{y}\sum_{i=1}^{n}\hat{y}_i + n\bar{y}^2 = \|\hat{\vec{y}}\|^2 - 2n\bar{y}^2 + n\bar{y}^2 = \|\hat{\vec{y}}\|^2 - n\bar{y}^2$$

$$\sum\hat{y} = \hat{\vec{y}}^T\vec{1}_n = (H\vec{y})^T\vec{1}_n = \vec{y}^T H^T\vec{1}_n = \vec{y}^T H\vec{1}_n = \vec{y}^T\vec{1}_n = \sum y_i = n\bar{y}$$

$\vec{1}_n \in colsp(X)$

$$= \sum_{j=0}^{p}\|proj_{\vec{e}_j}(\vec{y})\|^2 - n\bar{y}^2$$

What happens if we add a new feature to $X$   $\vec{x}_{new}$

$$X_{new} = [X | \vec{x}_{new}] \quad s.t. \quad X_{new} \text{ is full rank w/ rank} = (p+1)+1$$

What happens to SSR?

$$SSR_{new} = \sum_{j=0}^{p} \| proj_{\xi_j}(\vec{x}) \|^2 + \| proj_{\xi_{new}}(\vec{y}) \|^2 - n\bar{y}^2$$

ten predictors

$$> 0 \quad not = 0 \quad \text{otherwise } X^TX \text{ would not be full rank.}$$

$$\Rightarrow SSR_{new} = SSR + \| proj_{\xi_{new}}(\vec{y}) \|^2 > SSR$$

$$\Rightarrow R^2_{new} > R^2, \quad RMSE_{new} < RMSE$$

We can increase $R^2$ by just adding a vector to $X$. Even a totally made up vector!! lens...

What if we keep going.. keep adding columns until $p+1 = n$

$X$ is square $n \times n$ and full rank $\Rightarrow X$ is now invertible.

$$H = X(X^TX)^{-1}X^T = X X^{-1}(X^T)^{-1}X^T = I_n \quad why?$$

$$(AB)^{-1} = B^{-1}A^{-1}$$

Is $I_n$ a proj matrix? Yes... onto $\mathbb{R}^n$ the whole space.

Two properties

$I^T = I$ ✓

$I \cdot I = I$ ✓

eigenvalues of $I$? $\lambda_1 = \cdots = \lambda_n = 1$.

eigenvectors of $I$? Any basis for $\mathbb{R}^n$ will do...

$$\hat{y} = Hy = Iy = y \; ! \qquad \hat{y} = y \Rightarrow e = y - \hat{y} = \vec{0} \; !$$

$$\Rightarrow SSE = 0 \Longleftrightarrow RMSE = 0 \Longleftrightarrow R^2 = 100\%. \text{ Perfect fit!}$$

All you need to do is fill up $X$ with garbage random columns. Can this be real? NO!

It is called "in-sample" overfitting. (In-sample means within $D$). It is a core concept in this class.

How did this happen? $\vec{X}_{\cdot p+2}, \vec{X}_{\cdot p+3}, \dots \vec{X}_{\cdot n}$ are all reducing SSE by locating chance capitalization on different dimensions of the residual space.

These vectors are only chance-correlated with $y$ as they are not truly related to the $z$'s in $t(z)$. Thus, this chance correlation must disappear!

___

Why is this bad? Because this class is focused on using the model $g$ to make predictions on new data, $\hat{y}_* = g(\vec{x}_*)$. Overfitting is bad for future predictions. $\text{aka out-sample (OOS)}$

Bad for "generalization error" $\rightarrow$ how accurate is the model in the future?