

$X = QR$ decomposition has two steps 1) Gram-Schmidt algorithm which converts X into Q column-by-column and 2) reconstruction of the upper triangular change-of-basis matrix R . X has dimension $n \times k$ and columns x_1, \dots, x_k

In 1), we first a) create an orthogonal basis v_1, \dots, v_k and then b) normalize its component vectors into q_1, \dots, q_k .

a) $\vec{v}_1 = \vec{x}_1$

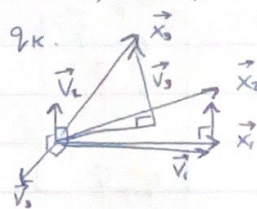
$\vec{v}_2 = \vec{x}_2 - \text{proj}_{\vec{v}_1}(\vec{x}_2)$

$\text{Span}\{\vec{x}_1, \vec{x}_2\} = \text{Span}\{\vec{v}_1, \vec{v}_2\}$ but $\vec{v}_1 \perp \vec{v}_2$

$\vec{v}_3 = \vec{x}_3 - \text{proj}_{[\vec{v}_1, \vec{v}_2]}(\vec{x}_3)$

\vdots

$\vec{v}_k = \vec{x}_k - \text{proj}_{[\vec{v}_1, \dots, \vec{v}_{k-1}]}(\vec{x}_k)$



b) $\vec{q}_1 = \frac{\vec{v}_1}{\|\vec{v}_1\|}$, $\vec{q}_2 = \frac{\vec{v}_2}{\|\vec{v}_2\|}$, \dots , $\vec{q}_k = \frac{\vec{v}_k}{\|\vec{v}_k\|}$

$\Rightarrow Q = [\vec{q}_1 | \dots | \vec{q}_k]$

2) $X = QR$

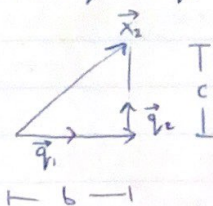
$[\vec{x}_1 | \vec{x}_2 | \dots | \vec{x}_k] = [\vec{q}_1 | \vec{q}_2 | \dots | \vec{q}_k] R$

$$R = \begin{bmatrix} a & b & d & & \\ 0 & c & e & & \\ 0 & 0 & f & \dots & \\ \vdots & \vdots & 0 & \dots & \\ 0 & 0 & 0 & \dots & \end{bmatrix}$$

$\vec{x}_1 = \|\vec{x}_1\| \vec{q}_1$

$\vec{x}_2 = b\vec{q}_1 + c\vec{q}_2 = H_1\vec{x}_1 + H_2\vec{x}_2$

$= \vec{q}_1 \underbrace{\vec{q}_1^T \vec{x}_2}_b + \vec{q}_2 \underbrace{\vec{q}_2^T \vec{x}_2}_c$



$\vec{x}_3 = d\vec{q}_1 + e\vec{q}_2 + f\vec{q}_3$
 $\vec{q}_1^T \cdot \vec{x}_3 \quad \vec{q}_2^T \cdot \vec{x}_3 \quad \vec{q}_3^T \cdot \vec{x}_3$

Sidebar: QR decomposition helps to speedup the OLS estimate computation in the following way:

$$\vec{b} = (X^T X)^{-1} X^T \vec{y}$$

↓

$$X^T X \vec{b} = X^T \vec{y} \Rightarrow (QR)^T QR \vec{b} = (QR)^T \vec{y} \Rightarrow R^T \overbrace{Q^T Q}^I R \vec{b} = R^T \overbrace{Q^T \vec{y}}^{\vec{z}} \\ \Rightarrow R^T R \vec{b} = R^T \vec{z} \Rightarrow R^{T^{-1}} R^T R \vec{b} = R^{T^{-1}} R^T \vec{z} \Rightarrow R \vec{b} = \vec{z}$$

e.g. $p+1=3$

by back-substitution

$$\begin{bmatrix} a & b & d \\ 0 & c & e \\ 0 & 0 & f \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix} \Rightarrow f b_2 = z_2 \Rightarrow b_2 = \frac{z_2}{f} \\ \Rightarrow c b_1 + e b_2 = z_1 \Rightarrow b_1 = \frac{1}{c} (z_1 - e \frac{z_2}{f}), \dots$$

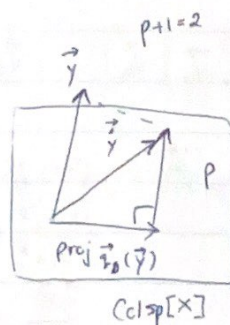
$R \quad \vec{b}$

$$SST = SSR + SSE \Rightarrow SSR \uparrow \Leftrightarrow SSE \downarrow$$

↑

fixed quantity only a function of \vec{y}

$$\begin{matrix} \updownarrow \\ R^2 \uparrow \Leftrightarrow RMSE \downarrow \end{matrix}$$



$$\hat{\vec{y}} = H\vec{y} = QQ^T \vec{y} = \sum_{j=0}^p \text{proj}_{\vec{q}_j}(\vec{y})$$

Pythag. thm

$$\|\hat{\vec{y}}\|^2 = \sum_{j=0}^p \|\text{proj}_{\vec{q}_j}(\vec{y})\|^2 = \underbrace{\|\text{proj}_{\vec{q}_0}(\vec{y})\|^2}_{\text{SST}} + \sum_{j=1}^p \|\text{proj}_{\vec{q}_j}(\vec{y})\|^2$$

$$= \|\text{proj}_{\vec{T}}(\vec{y})\|^2 = (H_0 \vec{y})^T (H_0 \vec{y}) = \vec{y}^T H_0 H_0 \vec{y} = (\vec{y} \vec{1}_n^T)(\vec{y} \vec{1}_n) = \vec{y}^T \vec{1} \vec{1}^T \vec{y} = n \bar{y}^2$$

$$SSR = (\vec{y} - \vec{y} \vec{1}^T)^T (\vec{y} - \vec{y} \vec{1}^T) = \vec{y}^T \vec{y} - \vec{y} \vec{1}^T \vec{y} - \vec{y} \vec{y}^T \vec{1} + \vec{y}^T \vec{1} \vec{1}^T \vec{y} \\ = \|\vec{y}\|^2 - 2 \bar{y} \vec{y}^T \vec{1} + n \bar{y}^2 = \|\vec{y}\|^2 - 2 n \bar{y}^2 + n \bar{y}^2 = \|\vec{y}\|^2 - n \bar{y}^2 = \sum_{j=1}^p \|\text{proj}_{\vec{q}_j}(\vec{y})\|^2$$

$(H\vec{y})^T \vec{1} = \vec{y}^T H^T \vec{1} = \vec{y}^T \vec{1} = n \bar{y}$

Pretend your friend gave you a new feature, i.e. a new \vec{x}_* . You want to now update your OLS model to use it.

$$X_* = [X \mid \vec{x}_*]$$

$$SSR_* = SSR + \underbrace{\|\text{proj}_{\vec{q}_*}(\vec{y})\|^2}_{>0} \Rightarrow SSR_* \geq SSR \Leftrightarrow SSE_* \leq SSE$$

Now your friend says "btw, I made up that vector... it's just a bunch of random nonsense". Any new column vector in X would have the ostensible effect of improving your model. If that new column is ~~not~~ independent of the true causal inputs to y (i.e. the z 's), we call this "overfitting".

Let's keep going. Your friend keeps supplying you with more and more garbage vectors. What happens when you have the same # of vectors $p+1 = n$?

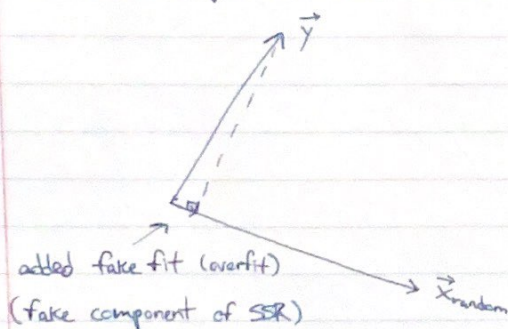
X_* will be $n \times n$ invertible so... $\text{colsp}[X_*] = \mathbb{R}^n$

$$H_* = X_*(X_*^T X_*)^{-1} X_*^T = X_* X_*^{-1} X_*^T X_* = I$$

$$\hat{y} = H_* y = y \Rightarrow \vec{e} = \vec{0}_n \Rightarrow \text{SSE} = 0 \Leftrightarrow R^2 = 1 \Leftrightarrow \text{RMSE} = 0$$

"perfect fit" or "maximal overfitting"

How did we get into this mess? Consider a random vec \vec{x}_{random} :



relative to n .

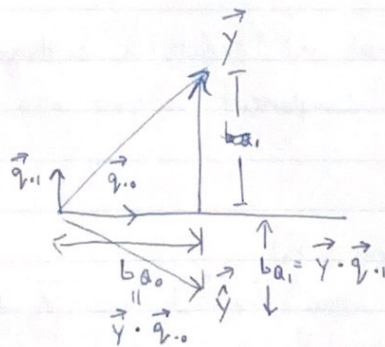
Overfitting becomes a problem with lots of features. If you have a small number of features, it's not too bad (i.e. it won't reduce your predictive accuracy).

We proved this in the context of OLS regression, but this is true in every other modeling context. Overfitting increases "generalization error" which is error on future predictions.

$$H = QR^T$$

$$\begin{aligned}\vec{b} &= (X^T X)^{-1} X^T \vec{y} = ((QR)^T (QR))^{-1} (QR)^T \vec{y} \\ &= (R^T \underbrace{Q^T Q}_I R)^{-1} R^T Q^T \vec{y} = (R^T R)^{-1} R^T Q^T \vec{y} = R^{-1} \underbrace{R^T R}_I Q^T \vec{y} \\ &= R^{-1} Q^T \vec{y}\end{aligned}$$

$$\vec{b}_Q = Q^T \vec{y} = \begin{bmatrix} \vec{q}_{\cdot 0}^T \vec{y} \\ \vec{q}_{\cdot 1}^T \vec{y} \\ \vdots \\ \vec{q}_{\cdot p}^T \vec{y} \end{bmatrix}$$



$$\text{Colsp}[\vec{q}_{\cdot 0}, \vec{q}_{\cdot 1}]$$

$$\text{Colsp}[\vec{1},$$