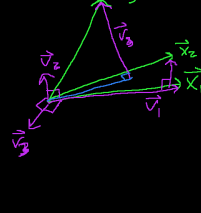X = QR decomposition has two steps (1) Gram-Shmidt algorithm which converts X into Q column-by-column and (2) reconstruction of the upper triangular change-of-basis matrix R. X has dimension n x K and columns x_1, ..., x_K.

In (1), we first (a) create a orthogonal basis v_1, ..., v_K and then (b) normalize its component vectors into q_1, ..., q_K.

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

$\vec{v}_2 := \vec{x}_2 - proj_{\vec{v}_1}(\vec{x}_2)$

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

$\vec{v}_3 := \vec{x}_3 - proj_{[\vec{v}_1 | \vec{v}_2]}(\vec{x}_3)$

$\vdots$

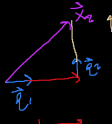$\vec{v}_K := \vec{x}_K - proj_{[\vec{v}_1 | \dots | \vec{v}_{K-1}]}(\vec{x}_K)$

(b) $\vec{q}_1 := \dfrac{\vec{v}_1}{\|\vec{v}_1\|}$ , $\vec{q}_2 := \dfrac{\vec{v}_2}{\|\vec{v}_2\|}$ , $\dots$ , $\vec{q}_K := \dfrac{\vec{v}_K}{\|\vec{v}_K\|}$

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

(2) $X = QR$   R

$[\vec{x}_1 | \dots | \vec{x}_K] = [\vec{q}_1 | \vec{q}_2 | \dots | \vec{q}_K]\begin{bmatrix} a & b & d & \cdots \\ 0 & c & e & \\ 0 & 0 & f & \\ \vdots & & & \ddots \\ 0 & 0 & & a \end{bmatrix}$

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

$\vec{x}_2 = b\vec{q}_1 + c\vec{q}_2 = H_1\vec{x}_2 + H_2\vec{x}_2 = \underbrace{\vec{q}_1 \vec{q}_1^T \vec{x}_2}_{b} + \underbrace{\vec{q}_2 \vec{q}_2^T \vec{x}_2}_{c}$

$\vec{x}_3 = d\vec{q}_1 + e\vec{q}_2 + f\vec{q}_3$

$\underset{\vec{q}_1 \cdot \vec{x}_3}{\|} \quad \underset{\vec{q}_2 \cdot \vec{x}_3}{\|} \quad \underset{\vec{q}_3 \cdot \vec{x}_3}{\|}$

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

$\vec{b} = (X^TX)^{-1}X^T\vec{y}$   very expensive operation

$\Downarrow$

$X^TX\vec{b} = X^T\vec{y} \Rightarrow (QR)^TQR\vec{b} = (QR)^T\vec{y} \Rightarrow R^T\overset{I}{\overbrace{Q^TQ}}R\vec{b} = R^TQ^T\overset{\vec{z}}{\overbrace{\vec{y}}}$

$\Rightarrow \overset{R^{-T}}{\overbrace{R^T}}R\vec{b} = R^T\overset{\vec{z}}{\overbrace{R^{-T}\vec{z}}} \Rightarrow R\vec{b} = \vec{z}$

e.g. $p+1 = 3$

$\begin{bmatrix} a & b & d \\ 0 & c & e \\ 0 & 0 & f \end{bmatrix}\begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix}$   by back-substition...

$\Rightarrow fb_2 = z_2 \Rightarrow b_2 = \dfrac{z_2}{f}$

$\Rightarrow cb_1 + eb_2 = z_1 \Rightarrow cb_1 + e\dfrac{z_2}{f} = z_1$

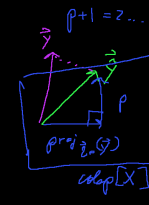$\Rightarrow b_1 = \dfrac{1}{c}(z_1 - e\dfrac{z_2}{f}), \dots e+c.$

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

$R^2\uparrow \quad \Leftrightarrow \quad RMSE\downarrow$

fixed quantity
only a function of   X = QR
the y-vector

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

by pythagoren thm.

$\|\vec{\hat{y}}\|^2 = \displaystyle\sum_{j=0}^{p}\|proj_{\vec{q}_j}(\vec{y})\|^2 = \|proj_{\vec{1}}(\vec{y})\|^2 + \displaystyle\sum_{j=1}^{p}\|proj_{\vec{q}_j}(\vec{y})\|^2$

$H_0 = \vec{1}\dfrac{(\vec{1}^T\vec{1})^{-1}}{\vec{1}^T} = \dfrac{1}{n}\begin{bmatrix} 1 & \cdots \\ \vdots & \end{bmatrix}$

$= \|proj_{\vec{1}}(\vec{y})\|^2 = (H_0\vec{y})^T(H_0\vec{y}) = (\vec{y}\,\vec{1}_n)^T(\vec{y}\,\vec{1}_n) = \vec{y}^2\,\vec{1}^T\vec{1} = n\vec{y}^2$

$SSR := (\vec{\hat{y}} - \bar{y}\vec{1})^T(\vec{\hat{y}} - \bar{y}\vec{1}) = \vec{\hat{y}}^T\vec{\hat{y}} - \bar{y}\vec{1}^T\vec{\hat{y}} - \bar{y}\vec{\hat{y}}^T\vec{1} + \bar{y}^2\vec{1}^T\vec{1}$

$= \|\vec{\hat{y}}\|^2 - 2\bar{y}\,\vec{\hat{y}}^T\vec{1} + n\bar{y}^2 = \|\vec{\hat{y}}\|^2 - 2n\bar{y}^2 + n\bar{y}^2 = \|\vec{\hat{y}}\|^2 - n\bar{y}^2 = \displaystyle\sum_{j=1}^{p}\|proj_{\vec{q}_j}\|^2$

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

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

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

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

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 vector is independent of the true causal inputs to y (i.e. the z's), we call this "overfitting".
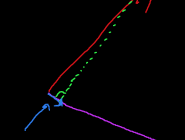
Let's keep going. Your friend keeps supplying you with more and more garbage vectors. What happens when you have the same number of vectors p+1 = n?

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

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

$\Rightarrow \vec{\hat{y}} = H_*\vec{y} = \vec{y} \Rightarrow \vec{e} = \vec{0}_n \Rightarrow SSE = 0 \Leftrightarrow R^2 = 1 \Leftrightarrow RMSE = 0$
"perfect fit" or "maximal overfitting"

How did we get into this mess? Consider a random vec x_random:

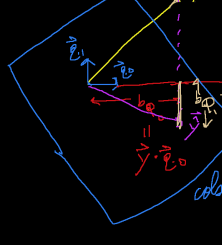added fake fit (over-fit)   $\vec{x}_{random}$
(fake component of SSR)

Overfitting becomes a problem with lots of features relative to n. If you have a small number of features relative to n, 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 modeling context. Overfitting increases "generalization error" which is error on future predictions.

$H = QR^T$

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

$= (R^TQ^TQR)^{-1}R^TQ^T\vec{y} = (R^TR)^{-1}R^TQ^T\vec{y}$

$= R^{-1}R^{-T}R^TQ^T\vec{y} = R^{-1}Q^T\vec{y}$

$\vec{b}_Q = Q^T\vec{y} = \begin{bmatrix} \vec{q}_0^T\vec{y} \\ \vec{q}_1^T\vec{y} \\ \vdots \\ \vec{q}_p^T\vec{y} \end{bmatrix}$

$colsp[\vec{q}_1, \vec{q}_2]$
$colsp[\vec{1}, \vec{x}_1]$