$X = QR$ decomposition has 2 steps (1) Gram-Shmidt algorithm which converts $X$ into $Q$ column-by-column and (2) reconstruction of the upper triangle change-of-basis matrix $R$. $X$ has dimension $n \times N$ and columns $x\_1, \ldots, x\_N$

In 1. we first (a) create a orthogonal basis $v\_1, \ldots, v\_N$ and then (b) normalize this component vectors into $q\_1, \ldots, q\_N$
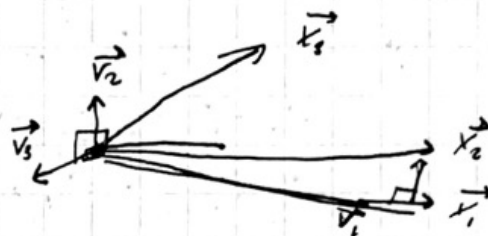
1. (a) $\vec{v_1} := \vec{x_1}$
$\vec{v_2} := \vec{x_2} - \text{proj}_{\vec{v_1}}(\vec{x_2})$

$\text{span}\{\vec{x_1}, \vec{x_2}\} = \text{span}\{\vec{v_1}, \vec{v_2}\}$ but $\vec{v_1} \perp \vec{v_2}$
$\vec{v_3} := \vec{x_3} - \text{proj}[\vec{v_1}|\vec{v_2}](\vec{x_3})$
$\vdots$
$\vec{v_N} := \vec{x_N} - \text{proj}[\vec{v_1}|\ldots|\vec{v_{N-1}}](\vec{x_N})$

1(b) $\vec{q_1} := \frac{\vec{v_1}}{\|\vec{v_1}\|}, \quad \vec{q_2} := \frac{\vec{v_2}}{\|\vec{v_2}\|}, \quad \ldots, \quad q_n := \frac{\vec{v_N}}{\|\vec{v_N}\|}$
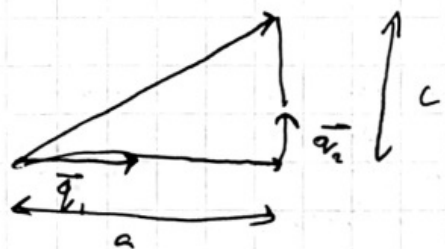
$\Rightarrow Q = [\vec{q_1}|\ldots|\vec{q_N}]$

② $X = QR$
$[\vec{x_1}|\vec{x_2})|\cdots|\vec{x_N}] = [\vec{q_1}|\vec{q_2}|\cdots|\vec{q_N}]$

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

$\vec{x_1} = \|\vec{x_1}\| \vec{q_1}$
$\vec{x_2} = b\vec{q_1} + c\vec{q_2} = H_1\vec{x_2} + H_2\vec{x_2} = \vec{q_1}\underbrace{\vec{q_1}^T\vec{x_2}}_{b} + \vec{q_2}\underbrace{\vec{q_2}^T\vec{x_2}}_{c}$



$\vec{x_3} = d\vec{q_1} + e\vec{q_2} + f\vec{q_3}$
$\quad \underset{\vec{q_1}\cdot\vec{x_3}}{\|} \quad \underset{\vec{q_2}\cdot\vec{x_3}}{\|} \quad \underset{\vec{q_3}\cdot\vec{x_3}}{\|}$

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

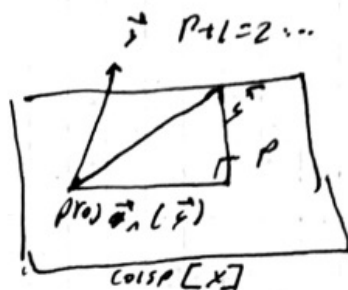$\vec{b} = (X^TX)^{-1}X^T\vec{y}$  very expensive operation.

$X^TX\vec{b} = \vec{y}^T\vec{y} \Rightarrow (QR)^TQR\vec{b} = (QR)^T\vec{y} \Rightarrow R^TQ^TQR\vec{b} = R^TQ^T\vec{y}$
$\Rightarrow R^{T-1}R^TR\vec{b} = R^{T-1}R^T \Rightarrow R\vec{b} = \vec{z}$  by back substitn.

e.g. $n+1 = 3$
$\Rightarrow \begin{bmatrix} a & b & d \\ 0 & c & e \\ 0 & 0 & f \end{bmatrix}\begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix}$

$\Rightarrow f b_2 = z_3 \Rightarrow b_2 = z_2/f$
$\Rightarrow cb_1 + eb_2 = z_1 \Rightarrow cb_1 + e\frac{z_2}{f} = z_1$
$\Rightarrow b_1 = \frac{1}{c}(z_1 - e\frac{z_2}{f}), \cdots$ etc

$\underset{R}{\ } \quad \underset{\vec{b}}{\ }$

$$SS\hat{y} = SSR + SSE \Rightarrow SSR\uparrow \Leftarrow\Rightarrow SSE\downarrow$$
$$R^2\uparrow \Leftarrow\Rightarrow RMSE\downarrow$$

Fixed once it's only
a function of the $\qquad X = QR$
$y$-vector



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

$$\|\vec{\hat{y}}\|^2 = \sum_{j=0}^{P}\|proj_{\vec{q}_j}(\vec{y})\|^2 = \|proj_{\vec{q}_0}(\vec{y})\|^2 + \sum_{j=1}^{P}\|proj_{\vec{q}_j}(\vec{y})\|^2$$

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

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

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

$$= \|\vec{\hat{y}}\|^2 - 2\bar{y}\,\vec{\hat{y}}^T\vec{1} + n\bar{y}^2 = \|\vec{\hat{y}}\|^2 - 2n\bar{y}^2 + n\bar{y}^2 = \|\vec{\hat{y}}\|^2 - n\bar{y}^2 =$$

$$\sum_{j=1}^{P}\|proj_{\vec{q}}(\vec{y})\|^2$$

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

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

$$X_f = [X \mid \vec{x}_m]$$

$$SSR_f = SSR + \|proj_{\vec{q}_f}(\vec{y})\|^2 \Rightarrow SSR_a \geq SSR \Leftrightarrow SSE_a \leq SSE$$
$$\underset{\geq 0}{\underbrace{\qquad}}$$

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 the 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 $p+1 = n$
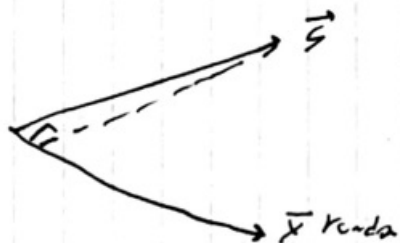
$$X_a \text{ will be } x \times n \text{ and invertible so } dim[X_a] = \mathbb{R}^n$$

$$H_a = X_a(X_a^TX_a)^{-1}X_a^T = X_a X_a^{-1} X_a^{T-1} X_a^T = I$$

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

"perfect fit" or "maximal overfitting"

How did we get into this mess? consider a random vec $\vec{x}$-random

$\vec{\hat{y}}$

added feature (over-fit)

(fake component of SSR)

$\bar{x}$ renda

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 modelling context, overfitting increases "generalization error" which is error on future predictions.

$$H = QR^T$$

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

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

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

$$\vec{ba} = Q^T\vec{y} = \begin{bmatrix} \vec{q}Q^T\vec{y} \\ q_2Q^T\vec{y} \\ \vec{q}_iQ^T\vec{y} \end{bmatrix}$$