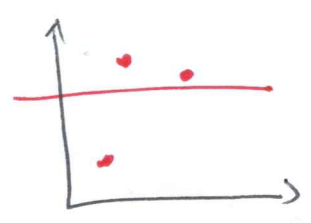


If $n = p+1 \Rightarrow X$ is square.

(no linear d.p.d.f.)

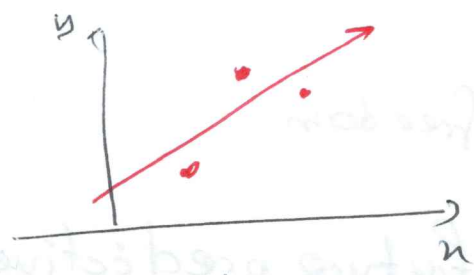
$$\Rightarrow H = X(X^T X)^{-1} X^T = X X^{-1} (X^T X^{-1})^T X^T = I$$

$$\vec{y} = H \vec{y} = I \vec{y} = \vec{y} \Rightarrow \vec{e} = \vec{y} - \vec{\hat{y}} = \vec{0} \Rightarrow SSE = 0 \Rightarrow R^2 = 100\%$$



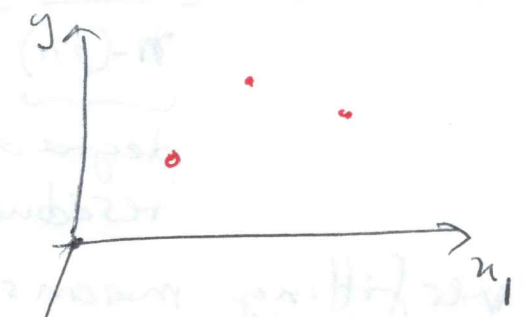
d.o.f = 1

"underfit"



d.o.f = 2

"fit well"

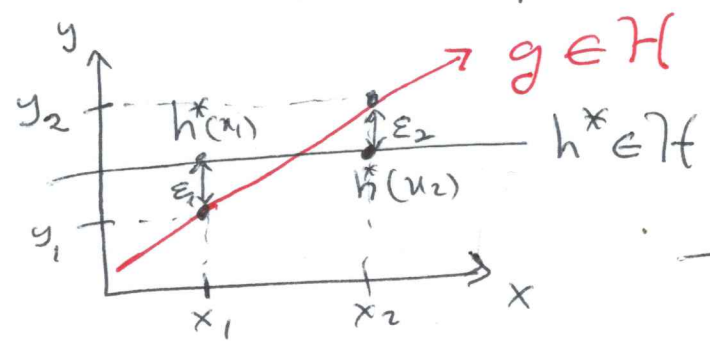


"overfit"

$$y = h^*(\vec{x}) + \epsilon$$

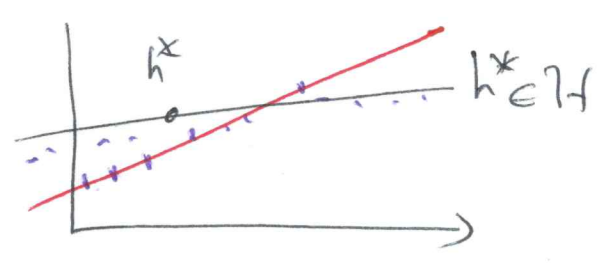
$h^* \in \mathcal{H}$ best model in candidate set given predict:

y has two components $h^*(\vec{x})$, ϵ (sympo, noise)



overfit

The overfix model g doesn't generalize (doesn't fit random well)



$$y = g(\vec{x}) + e$$

$$= g(\vec{x}) + \underbrace{(h^* - g)}_{\text{estimat}^{\circ} \text{ error (this includes ordinary error)}} + \varepsilon$$

$n \rightarrow \infty$ \exists "estimation error" $\rightarrow 0$

and p constant.

Note: $MSE = \frac{1}{\underbrace{n - (p+1)}_{\text{degree of freedom residual}}} SSE$

Overfitting means future predictive performance
 You only see overfitting iff you never see the data before
 g and h^* are the same based on n observations

