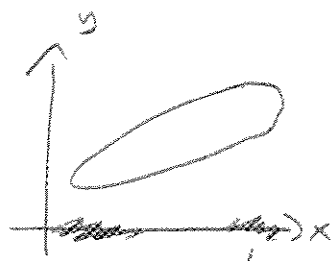


Least Squares estimates

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_j (x_j - \bar{x})^2} = \sum_i k_i \cdot y_i, \quad k_i = \frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2}$$

linear comb  
of all y-values



Extremes in x dominate  
the fit

Properties

$$E(\hat{\beta}_1) = \beta_1$$

unbiased - so average  
across multiple data sets

we get true model

but we only see one

data in practise - how much

can we be off in a typical

case?

Variance

$$V(\hat{\beta}_1) = \frac{\sigma^2}{\sum_j (x_j - \bar{x})^2}$$

$$\hat{\beta}_1 = \beta_1 \text{ (give or take)} \underbrace{\sqrt{V(\hat{\beta}_1)}}_{\text{captures uncertainty}} \text{ essentially}$$

Captures uncertainty

### 3 sources of uncertainty

①  $V(\hat{\beta}_1) \uparrow$  as  $\sigma^2 \uparrow$  - i.e. more noise in the data makes estimation more difficult

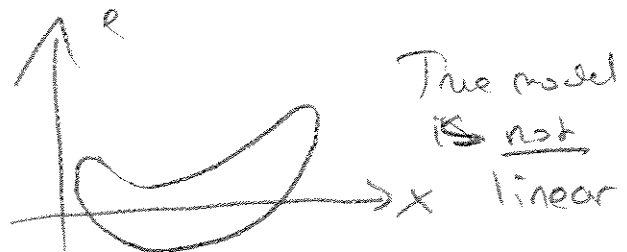
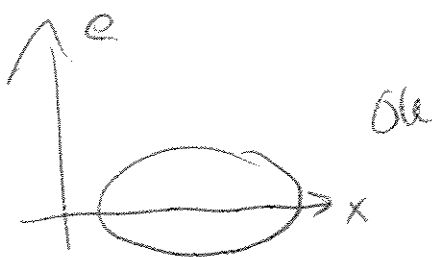
②  $V(\hat{\beta}_1) \downarrow$  as  $\frac{1}{n} \uparrow$  - i.e. more data makes estimation easier

③  $V(\hat{\beta}_1) \downarrow$  as  $V(x) \uparrow$  - i.e. more spread in  $x$  (as captured by  $\sum(x_i - \bar{x})^2$ ) makes estimation easier.

### Orthogonality properties

•  $\sum_i x_i e_i = 0 \Rightarrow$  i.e. residuals are orthogonal to  $x$   
 $e$  and  $x$  uncorrelated

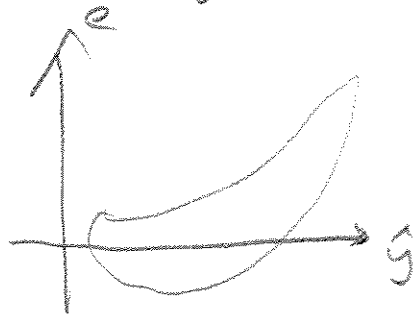
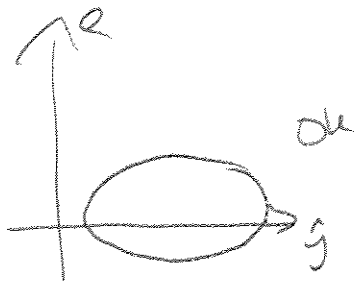
$\Rightarrow$  we used up all the linear information in  $x$  about  $y$



$$\sum_{i=1}^n \hat{y}_i \cdot e_i = 0$$

i.e. residuals are orthogonal to fitted values

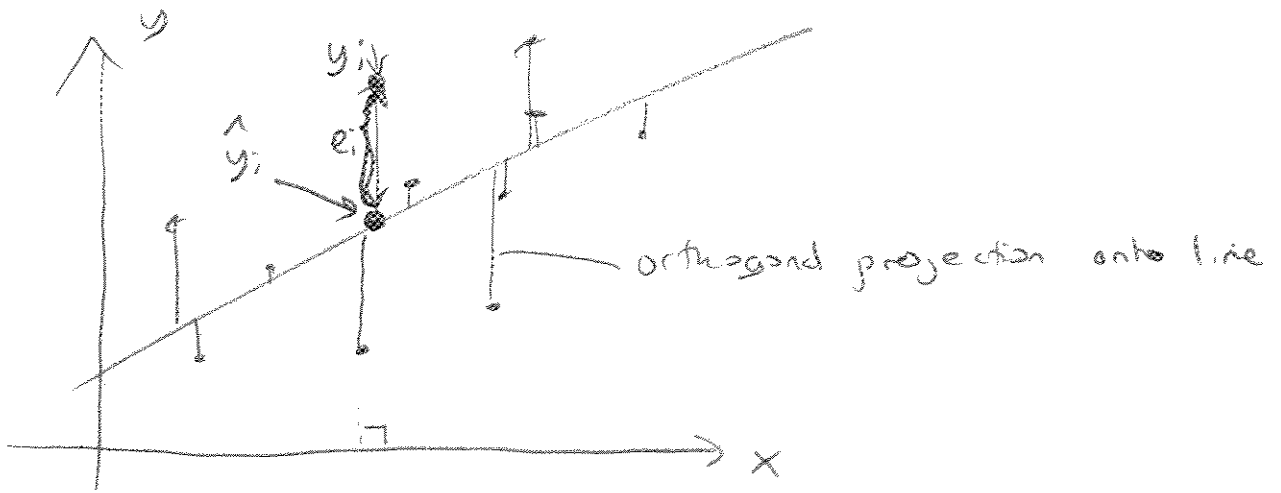
$\Rightarrow e$  and  $\hat{y}$  are uncorrelated



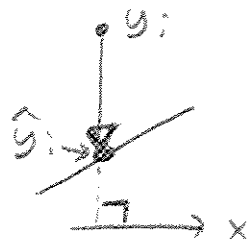
model insufficient.

follows from

$\hat{y}$  = a function of  $x$  and  $x$  is uncorrelated with  $e$ .



So  $\hat{y}_i = y_i$  projected onto the line orthogonal to  $x$



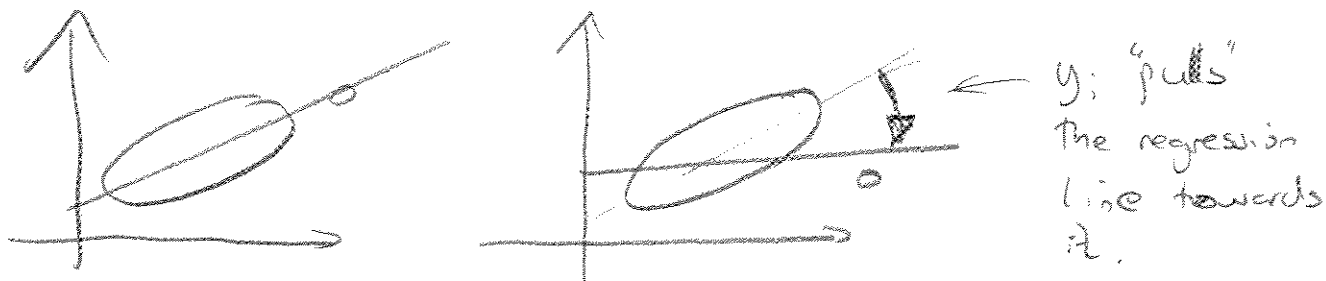


if  $h_{ii} \gg h_{ij}$ , i.e.  $y_i$  dominates the fit at  $x_i$  (5)

- we call  $x_i$  a point of high leverage



Problem : if  $x_i =$  high leverage and corresponding  $y_i$  is extreme  $\Rightarrow$  poor fit



Where does the leverage show up in properties of fitted values?

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \implies E(\hat{y}_i) = E(y_i) = \beta_0 + \beta_1 x_i \text{ since } \hat{\beta}_0 \text{ and } \hat{\beta}_1 \text{ are unbiased}$$



$$V(\hat{y}_i) = V(\hat{\beta}_0 + \hat{\beta}_1 x_i) = V(\bar{y} + \hat{\beta}_1 (x_i - \bar{x}))$$

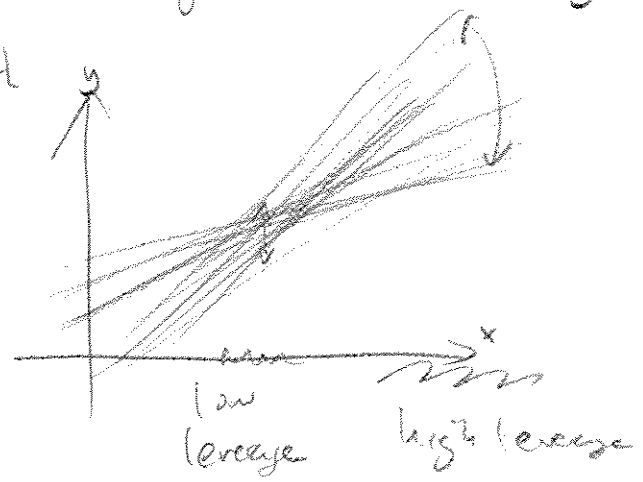
$$= V(\bar{y}) + (x_i - \bar{x})^2 V(\hat{\beta}_1) \text{ because } \bar{y} \text{ and } \hat{\beta}_1 \text{ are uncorrelated}$$

$$\left[ \begin{aligned} \text{Cov}(\bar{y}, \hat{\beta}_1) &= \text{Cov}\left(\frac{1}{n} \sum y_j, \sum k_j \cdot y_j\right) \\ &= \sum_j \frac{k_j}{n} \text{Cov}(y_j, y_j) \text{ [since } y_i \text{ and } y_j \text{ are uncorrelated]} \\ &= \sum_j \frac{k_j}{n} V(y_j) = \frac{\sigma^2}{n} \sum_j k_j = \frac{\sigma^2}{n} \cdot 0 = 0 \end{aligned} \right]$$

$$= \frac{\sigma^2}{n} + \frac{\sigma^2}{\sum_j (x_j - \bar{x})^2} (x_i - \bar{x})^2 = \sigma^2 \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_j (x_j - \bar{x})^2} \right) = h_{ii} \cdot \sigma^2$$

↑  
leverage

So, at points of high leverage, (lots of) uncertainty in estimate  $\hat{y}_i$ , since small fluctuations in  $y_i$  can really alter the fit



What about the residuals?

⑦

$$e_i = y_i - \hat{y}_i, \quad E(e_i) = 0$$

$$V(e_i) = \text{Cov}(e_i, e_i) = \text{Cov}(y_i - \hat{y}_i, y_i - \hat{y}_i)$$

$$= \text{Cov}(y_i, y_i) + \text{Cov}(\hat{y}_i, \hat{y}_i) - 2 \text{Cov}(y_i, \hat{y}_i)$$

$$= \sigma^2 + \sigma^2 h_{ii} - 2 \text{Cov}(\hat{y}_i + e_i, \hat{y}_i)$$

$$= \sigma^2 + \sigma^2 h_{ii} - 2 \sigma^2 h_{ii} - 2 \text{Cov}(e_i, \hat{y}_i)$$

$$= \underline{\underline{\sigma^2(1-h_{ii})}}$$

" 0 since  $e_i$  uncorrelated with  $\hat{y}_i$ .

So, residual variance smaller at points of high leverage since whatever  $y_i$  is it pulls the regression line towards it — controls the residual magnitude.

Note also  $\text{Cov}(e_i, e_j) = -\sigma^2 h_{ij}$  so residuals are correlated

True errors,  $\epsilon_i$

$$E(\epsilon_i) = 0$$

$$V(\epsilon_i) = \sigma^2 \text{ constant}$$

$$\text{Cov}(\epsilon_i, \epsilon_j) = 0 \text{ uncorrelated}$$

Residuals from fit,  $e_i$

$$E(e_i) = 0$$

$$V(e_i) = \sigma^2(1-h_{ii}) \text{ not constant}$$

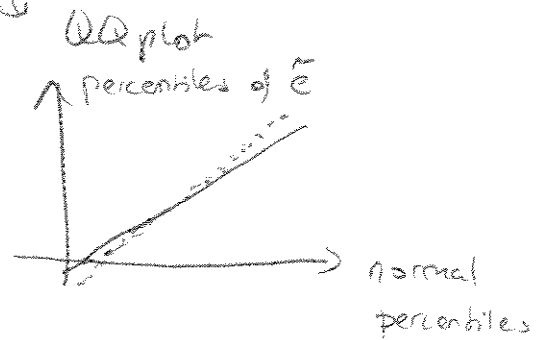
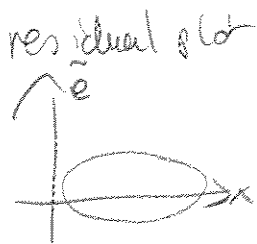
$$\text{Cov}(e_i, e_j) = -\sigma^2 h_{ij} \text{ correlated.}$$

⇒ Careful when looking at residual plots!

⑧

Since  $V(e_i) = \sigma^2(1-h_{ii})$ , can only directly compare residuals if we adjust for the non-constant variance

→ Standardized residuals  $\tilde{e}_i = \frac{e_i}{\sqrt{1-h_{ii}}}$ ,  $V(\tilde{e}_i) = \sigma^2$   
Constant





# The nuisance parameter, $\sigma^2$

9

$V(\hat{\beta}_1), V(\hat{y}_i), V(e_i)$  all involve  $\sigma^2$  — so we can't really do inference without knowing it.

Use data to estimate both  $\beta$  and  $\sigma^2$

$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ ,  $E(\varepsilon_i) = 0$ ,  $V(\varepsilon_i) = \sigma^2$  assumed model

So  $\sigma^2 =$  variance of true errors  $\varepsilon$ .

We have the LS fit  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , and the residuals  $e$ .

Can we use  $e$  to estimate  $\sigma^2$ ?

If we knew  $\varepsilon_1, \dots, \varepsilon_n$  (true errors)

$$(a) \hat{\sigma}^2 = V(\varepsilon) = \frac{1}{n-1} \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2$$

Here, we use Mean Squared Error, MSE

$$(b) \text{MSE} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2} = \frac{\text{RSS}}{n-2}$$

Residual Sum of Squares  
↓

Why  $n-2$  not  $n-1$ ?

- In (a), the  $\varepsilon_i - \bar{\varepsilon}$  share one parameter  $\bar{\varepsilon}$ , i.e. we constrain the sum of  $\varepsilon_i - \bar{\varepsilon}$  to 0  $\Rightarrow$   $n-1$  degrees of freedom
- In (b),  $y_i - \hat{y}_i$  share two parameters,  $\hat{\beta}_0$  and  $\hat{\beta}_1$ ,  $\Rightarrow$   $n-2$  degrees of freedom

MSE  $\Rightarrow$  unbiased estimate of  $\sigma^2$

Key facts used:  
 $V(Z) = E(Z^2) - (E(Z))^2$   
 $\text{cov}(Z, W) = E(ZW) - E(Z)E(W)$

$$E(RSS) = E\left(\sum_{i=1}^n (y_i - \hat{y}_i)^2\right) = E\left(\sum e_i^2\right)$$

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

$$= \sum_{i=1}^n \left( V(y_i) + (E(y_i))^2 + V(\hat{y}_i) + (E(\hat{y}_i))^2 - 2E(\hat{y}_i + e_i) \hat{y}_i \right)$$

$$= \sum_{i=1}^n V(y_i) + V(\hat{y}_i) + 2(E(y_i))^2 - 2\text{cov}(\hat{y}_i + e_i, \hat{y}_i) - 2E(y_i)E(\hat{y}_i)$$

$$= \sum_{i=1}^n V(y_i) + V(\hat{y}_i) + 2(E(y_i))^2 - 2(E(y_i))^2 - 2\text{cov}(\hat{y}_i, \hat{y}_i) - 2\text{cov}(e_i, \hat{y}_i)$$

$\begin{matrix} \parallel & & \parallel \\ V(\hat{y}_i) & & 0 \end{matrix}$

$$= \sum_{i=1}^n V(y_i) - V(\hat{y}_i)$$

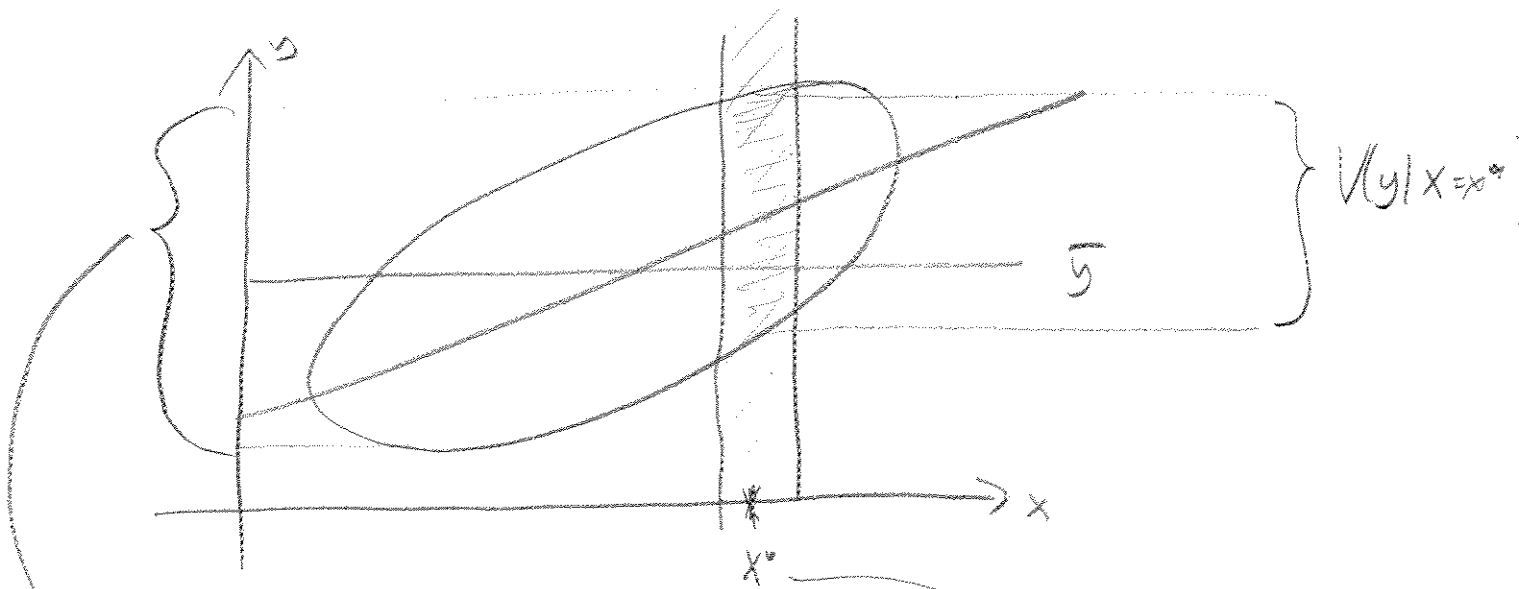
$$= n \cdot \sigma^2 - \sigma^2 \sum_{i=1}^n h_{ii} = \sigma^2 \left( n - \sum_{i=1}^n h_{ii} \right) = \sigma^2 \left( n - \sum_{i=1}^n \left( \frac{1 + (x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) \right)$$

$$= \sigma^2 (n-2)$$

$$\Rightarrow \boxed{E(MSE) = \sigma^2} \quad \#$$

# Variance Decomposition

(11)



Spread among  $y$ -values around  $\bar{y} \Rightarrow$  Marginal distribution for  $y$   
 $\Rightarrow$  marginal variance  $V(y)$

In neighborhood of  $x^*$  - spread among  $y$  in strip around regression value  $\hat{\beta}_0 + \hat{\beta}_1 x^*$   
 $\sim V(y|x=x^*)$   
= Conditional variance

If  $y$  and  $x$  are unrelated  $\Rightarrow V(y) \approx V(y|x)$   
- doesn't help to condition

If  $y = \beta_0 + \beta_1 x$  exactly  $\Rightarrow V(y|x) = 0$  - no variability left!

Note, in LS we assume  $V(\epsilon) = \sigma^2$  constant, or equivalently, that  $V(y|x=x^*)$  independent of actual  $x^*$  value

Summarize our fit  $\rightarrow$  are we close to

- unrelated  $x$  and  $y$

or

- perfect fit  $y = \beta_0 + \beta_1 x$  exactly?

Two quantities to compare:

$$\textcircled{1} \text{RSS} = \text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Spread among  $y_i$  around regression line

Residual sum of squares

Error sum of squares

two different names

$$\textcircled{2} \text{SST} = \sum_{i=1}^n (y_i - \bar{y})^2$$

Spread among  $y_i$  around mean of  $y$ .

Total sum of squares

Note, since we fit the regression line to data to minimize RSS, RSS is always smaller than SST

# The R-squared

13

$$R^2 = \frac{SS_T - RSS}{SS_T} = \frac{\text{reduction in spread of } y}{\text{total spread in } y}$$

= "% of variability in y explained by the regression"

