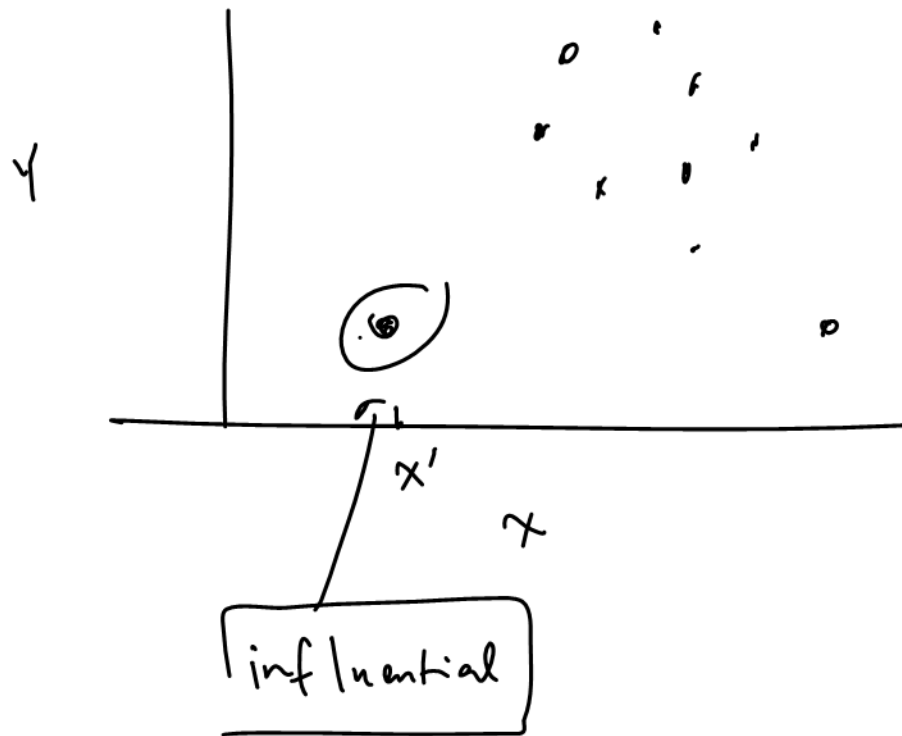


Lecture 09 Oct 11

Outlier & Influential Observations



$$\hat{\beta}_1(\text{complete}) \quad \text{vs} \quad \hat{\beta}_1(\text{without } (x', y'))$$

The estimator $\hat{\beta}_1$:

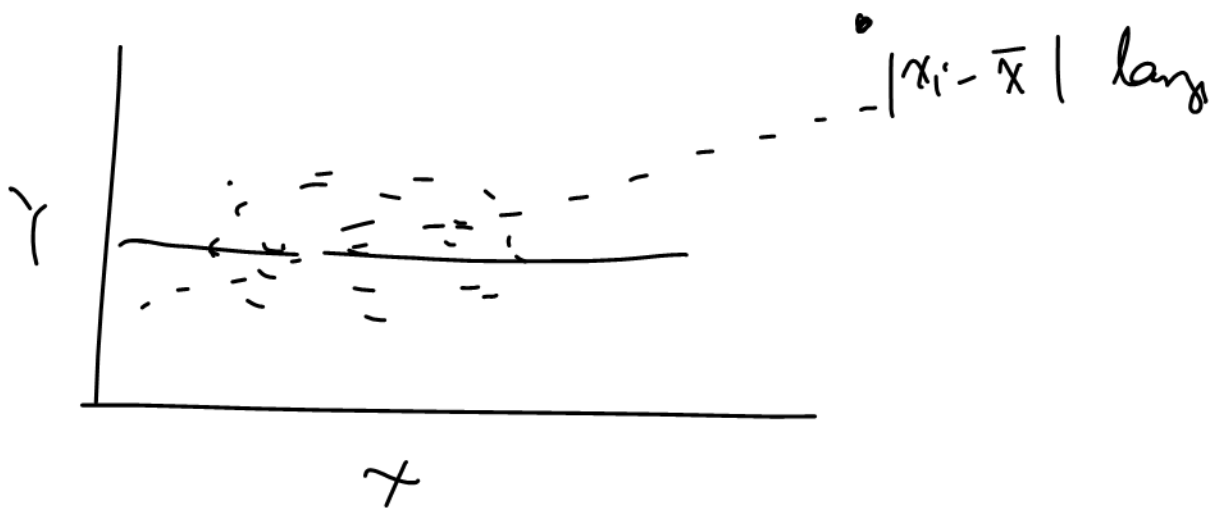
$$\hat{\beta}_1 = \sum w_i y_i$$

$$= \sum \left[\frac{(x_i - \bar{x})}{\sum (x_j - \bar{x})^2} \right] \cdot y_i$$

For values $\{x_i\}$ where x_i close to \bar{x}

$$|x_i - \bar{x}| < \delta \Rightarrow w_i \text{ small}$$

When $x_i - \bar{x}$ large (in magnitude)
then y_i has more weight in
estimating β_1 .



Let $(x_i, y_i) = \underline{z}_i$ be a candidate outlier/
influential observation.

$$\text{Let } \mathcal{D} = \{ (x_n, y_n), n=1, \dots, N \}$$

$$\mathcal{D}_{(-i)} = \mathcal{D} - \{ (x_i, y_i) \}$$

$$| \hat{\beta}_1(\mathcal{D}) - \hat{\beta}_1(\mathcal{D}_{(-i)}) |$$

Multi-collinearity

$$X = \begin{pmatrix} \mathbb{1} & x_1 \\ & x_2 \\ & \vdots \\ & x_n \end{pmatrix}$$

Suppose that

$$x_i = 10$$

$$\begin{pmatrix} \mathbb{1} & 10 \\ & \vdots \\ & 10 \end{pmatrix}$$

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

$$\hat{\beta} = (X'X)^{-1} X'Y$$

$$(X'X) = \begin{pmatrix} n & 10n \\ 10n & n(10^2) \end{pmatrix}$$

$$(X'X)^{-1} = \frac{1}{\det(X'X)} \begin{pmatrix} & \\ & \end{pmatrix}$$

$$\det(X'X) = n^2(10^2) - (10n)^2 = 0.$$

Ridge Regression

If $(X'X)$ is near singular, i.e., $\det(X'X) \approx 0$

$$\hat{\beta} = (X'X)^{-1} X'Y$$

$$\hat{\beta}_R = [(X'X) + \lambda I]^{-1} X'Y, \quad \lambda > 0$$

Diagnosis

Model $Y_i | x_i \stackrel{\text{indep}}{\sim} \mathcal{N}(\mu(x_i), \sigma^2)$

$$\mu(x_i) = \beta_0 + \beta_1 x_i \leftarrow$$

(1) Is linearity sufficient?

Is the linear structure reasonable?

(2) Normality

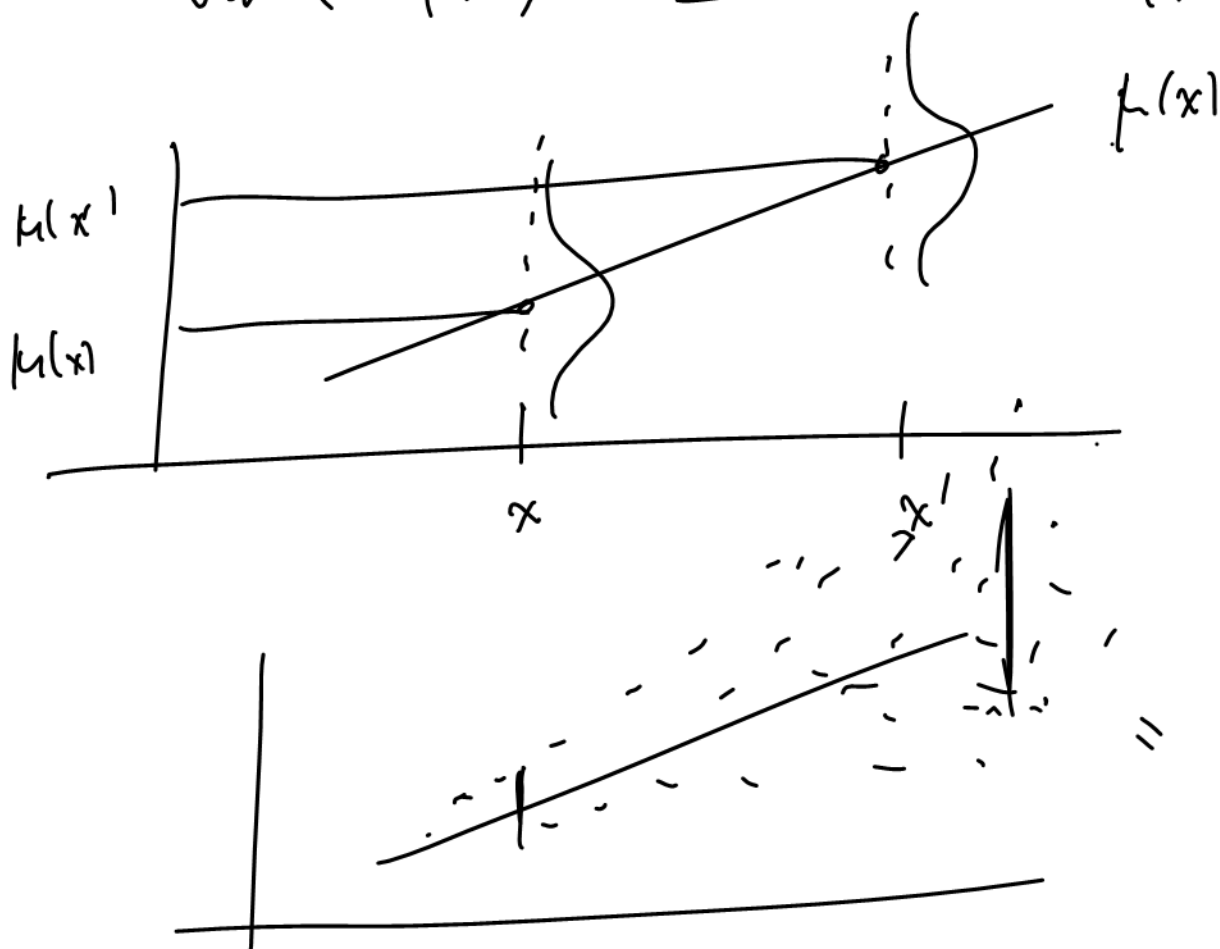
(3) Equal variance across all x values
Homoscedasticity

Heteros. - -

(4) Independence

According to the model: let x and x' be two values of the indep variable. Then:

$$\text{Var}(Y|x) = \text{Var}(Y|x')$$

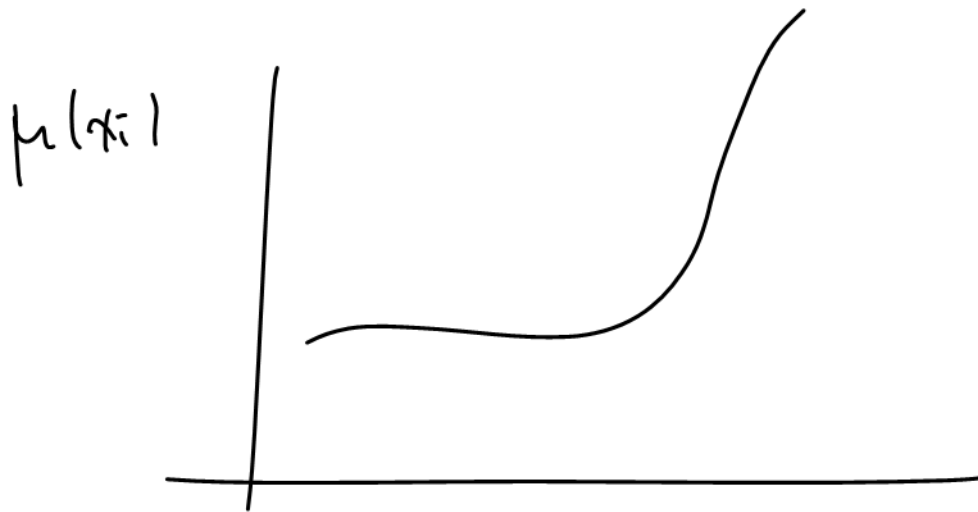
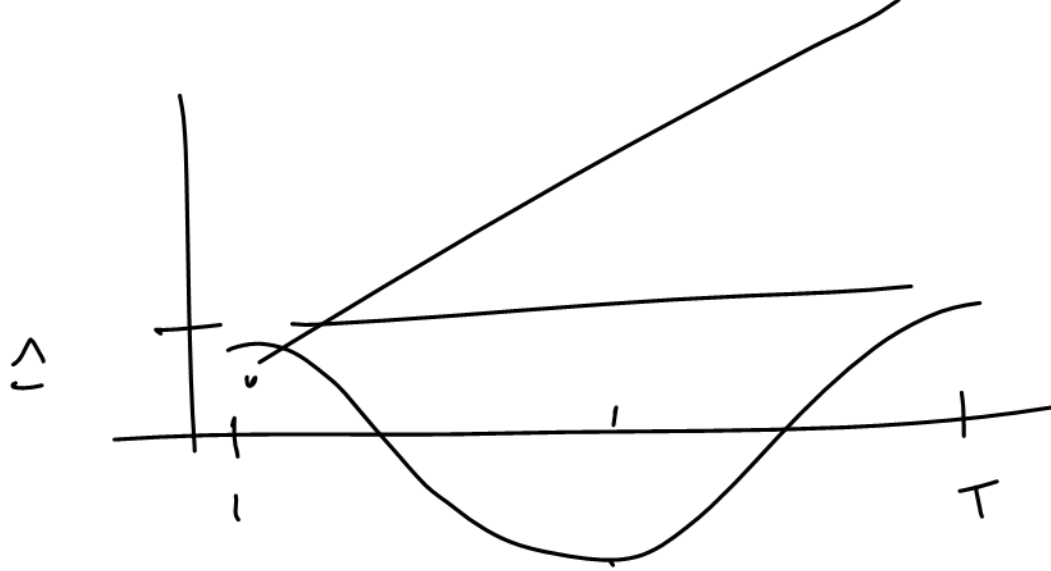


Diagnostics on the mean structure

Suppose that this is the true model:

$$Y_i|x_i \sim N(\mu(x_i), \sigma^2) \quad x_i = 1, 2, \dots, T$$

$$\mu(x_i) = \beta_0 + \beta_1 x_i + \alpha_1 \cos\left(2\pi \frac{x_i}{T}\right)$$



Fit this model :

$$Y_i | x_i \sim N(\mu(x_i), \sigma^2)$$

$$\mu(x_i) = \beta_0 + \beta_1 x_i$$

$$\Leftrightarrow \begin{aligned} Y_i &= (\beta_0 + \beta_1 x_i) + \varepsilon_i, \\ \varepsilon_i &\stackrel{\text{iid}}{\sim} N(0, \sigma^2) \end{aligned}$$

$$\varepsilon_i = Y_i - \mu(x_i) \quad \leftarrow \begin{array}{l} \text{not} \\ \text{observed} \end{array}$$

Residuals

$$\begin{aligned} R_i &= Y_i - \hat{\mu}(x_i) \\ &= Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \end{aligned} \quad \left. \vphantom{\begin{aligned} R_i &= Y_i - \hat{\mu}(x_i) \\ &= Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \end{aligned}} \right\} \begin{array}{l} \text{observed} \\ \text{from the} \\ \text{data} \end{array}$$

$$\begin{aligned} R_i &= \overbrace{\mu(x_i)}^{\text{True}} + \varepsilon_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \\ &= \underbrace{\beta_0 + \beta_1 x_i + \alpha_1 \cos\left(2\pi \frac{x_i}{T}\right)}_{Y_i} + \varepsilon_i \\ &\quad - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \end{aligned}$$

$$\begin{aligned} &= (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)x_i + \underline{\underline{\alpha_1 \cos\left(2\pi \frac{x_i}{T}\right)}} \\ &\quad + \varepsilon_i \end{aligned}$$

$$\approx \underline{\underline{\alpha_1 \cos\left(2\pi \frac{x_i}{T}\right)}} + \varepsilon_i$$

Plot of $\{R_i\}$



ANALYSIS OF
COVARIANCE