

this time: regression
next time: ANOVA

read: LN pp. L-269-289

AM57
18 Jul
2016
①

today: LN pp. L-247

homework 4 due Fri 22 Jul in class

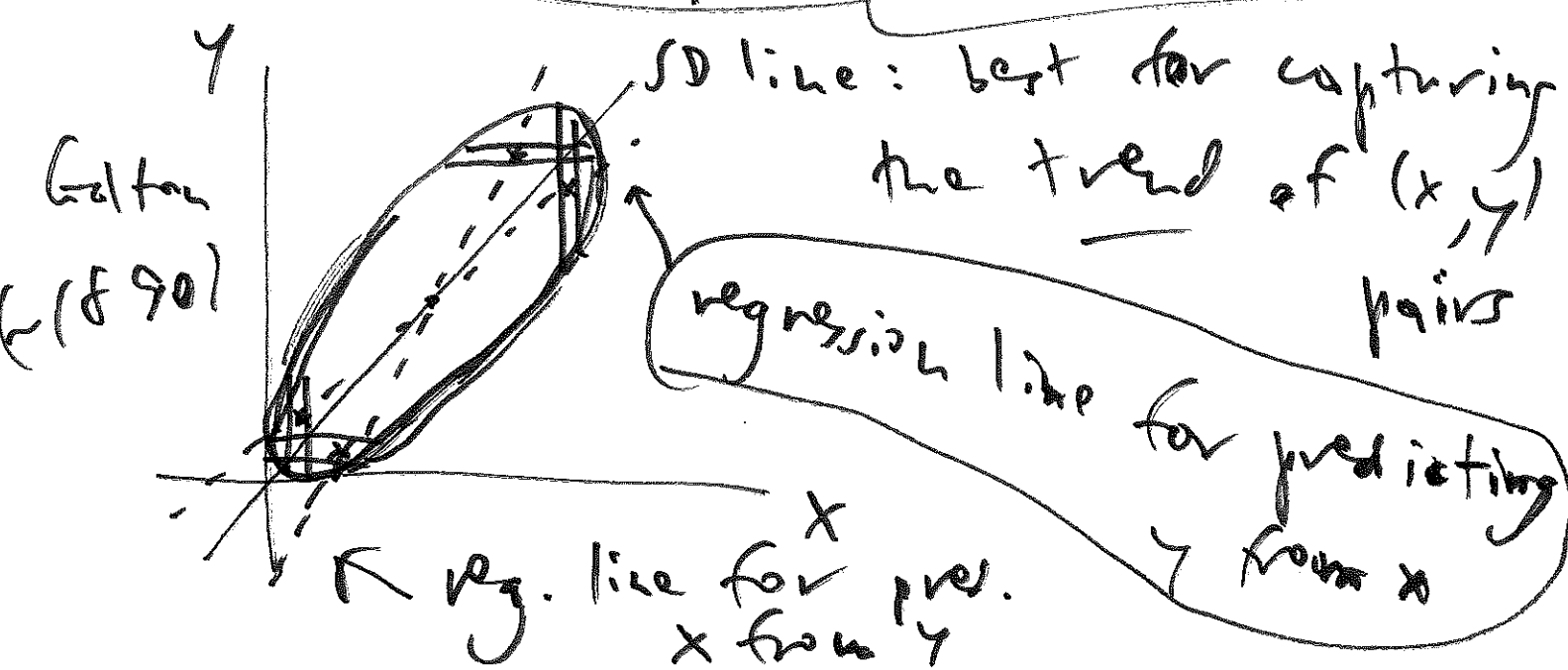
take-home final due Mon ~~22~~ 25 Jul noon

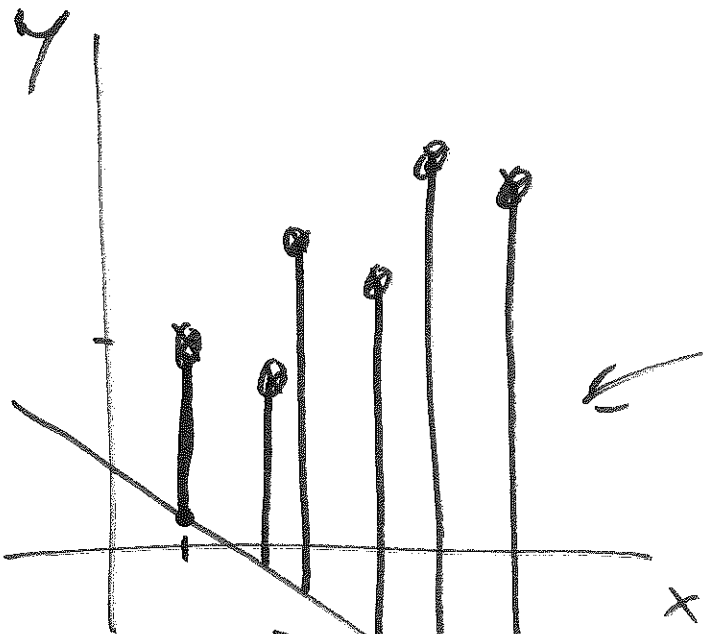
extra office hours: Sat 23 Jul noon-1pm

Mon 25 Jul 9A-10A Sun 24 Jul noon-1pm

take elevator to 3rd floor parking: door ripped open Sat 23 Jul 2-3p

Sun 24 Jul 2-3p

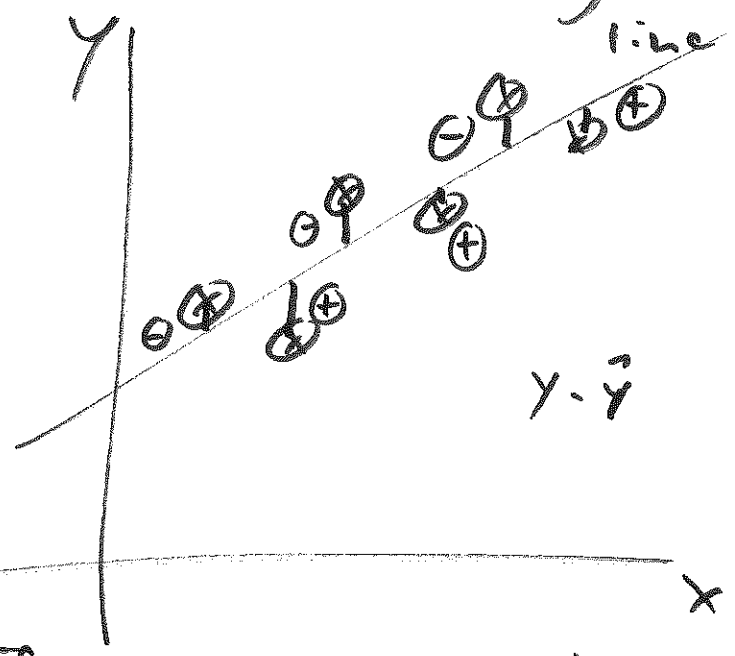




Goal: how predict y from x ?
 (residuals)

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$
 make predictions for y

a bad line



$y - \hat{y}$

find $(\hat{\beta}_0, \hat{\beta}_1)$ to minimize

$$\sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

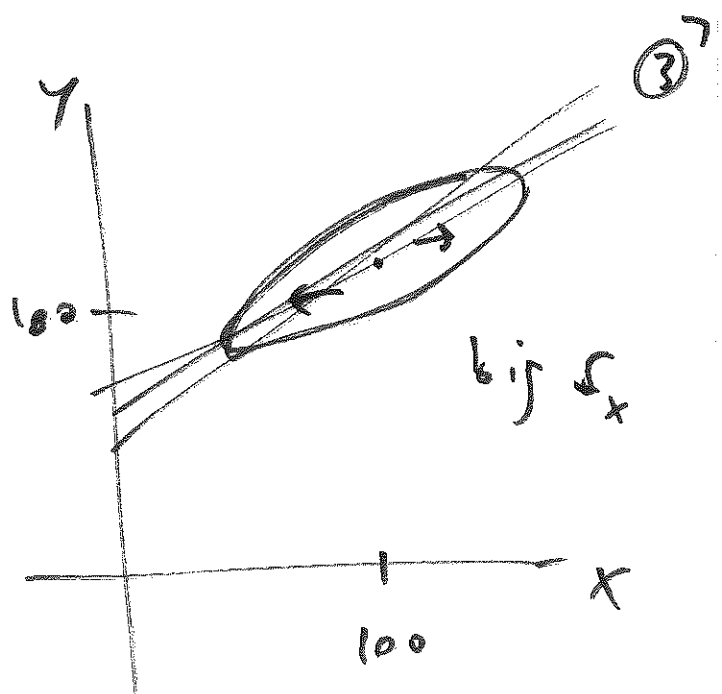
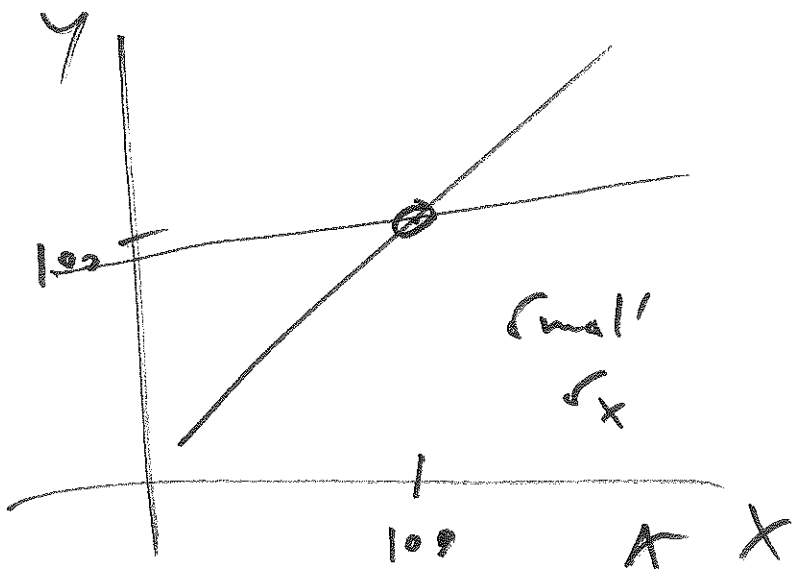
\hat{y}_i

$$\begin{pmatrix} y_1 & x_1 \\ \vdots & \vdots \\ y_n & x_n \end{pmatrix}$$

least sq have
line

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

math fact: regression line = least squares line



harder to estimate β_1 here

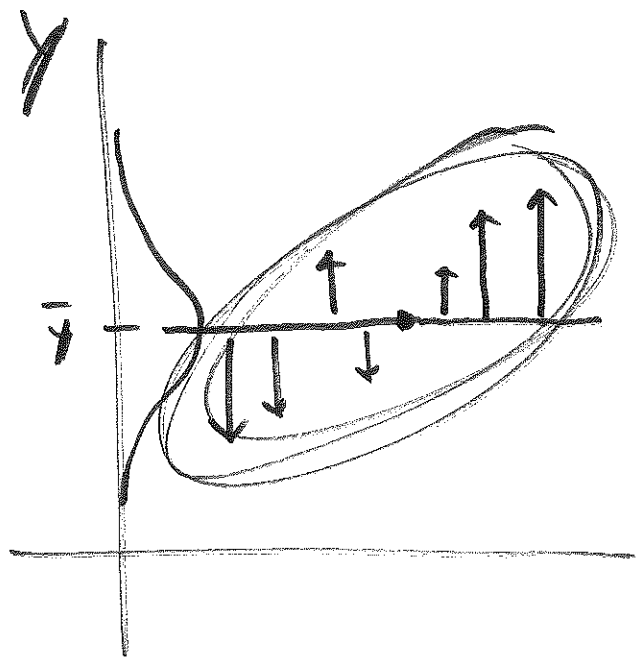
to judge whether

a sample s_b for $\beta_1 = \frac{s_y}{s_x} r$ is large in practical terms, use exactly the same line of reasoning as with the sample correlation r

(10.00)

(10.27)

how useful
is the regression?



$$\begin{pmatrix} y_1 & x_1 \\ \vdots & \vdots \\ y_n & x_n \end{pmatrix}$$

predictive task
① ignore x or don't measure it,

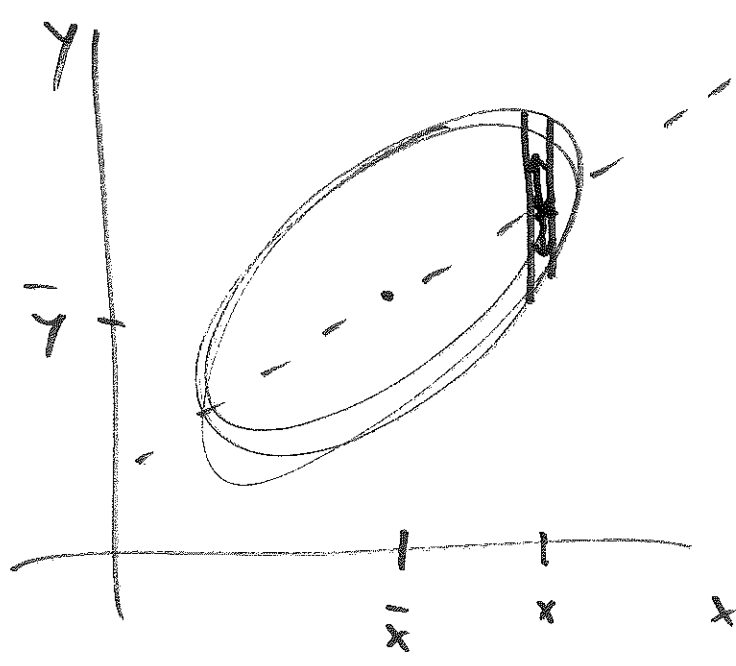
predict

y :

$$\hat{y}_{no\ x} = \bar{y}$$

with $SE(\hat{y}_{no\ x}) = s_y$

predictive task ②: ⑤

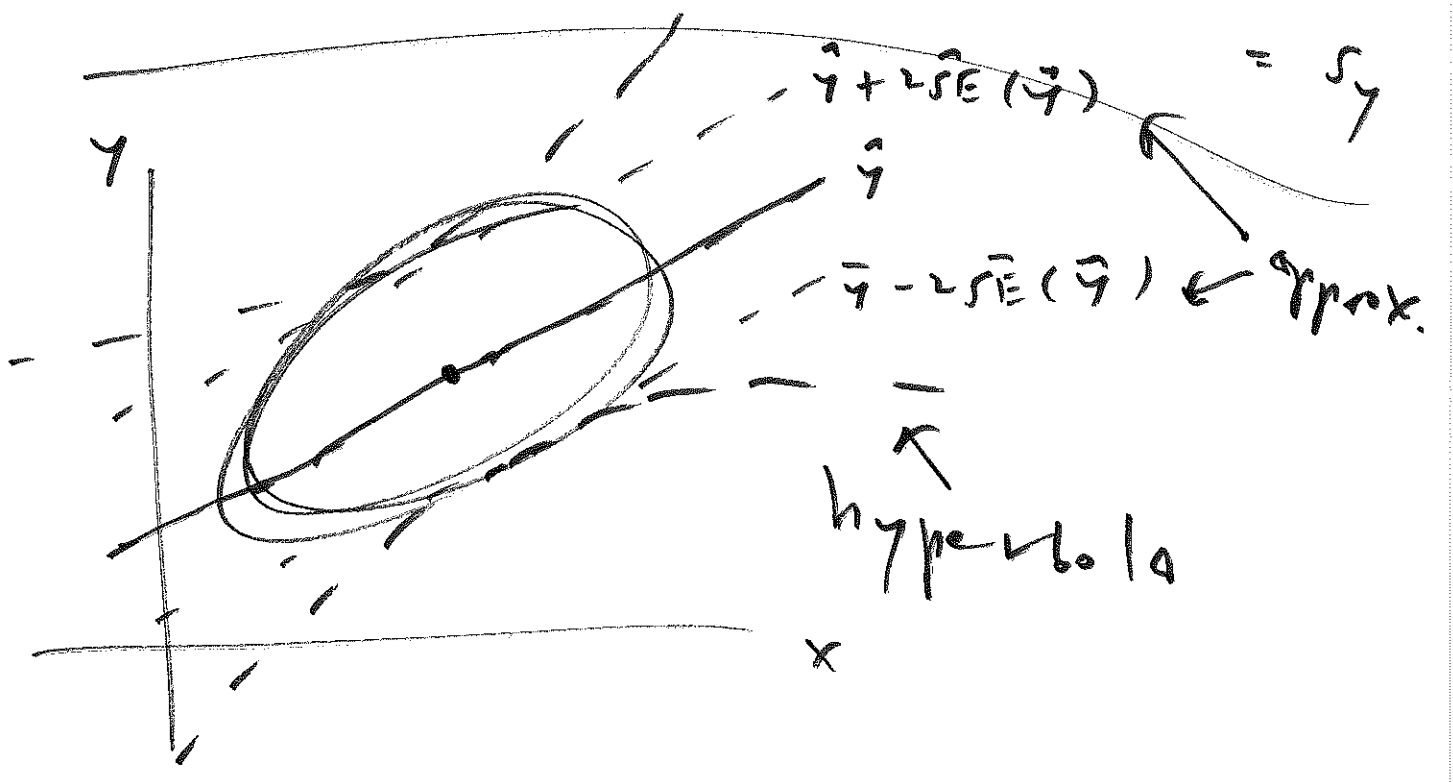


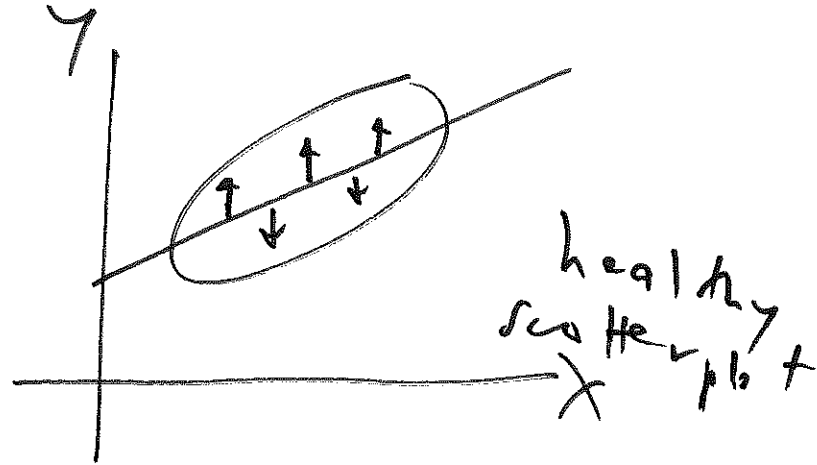
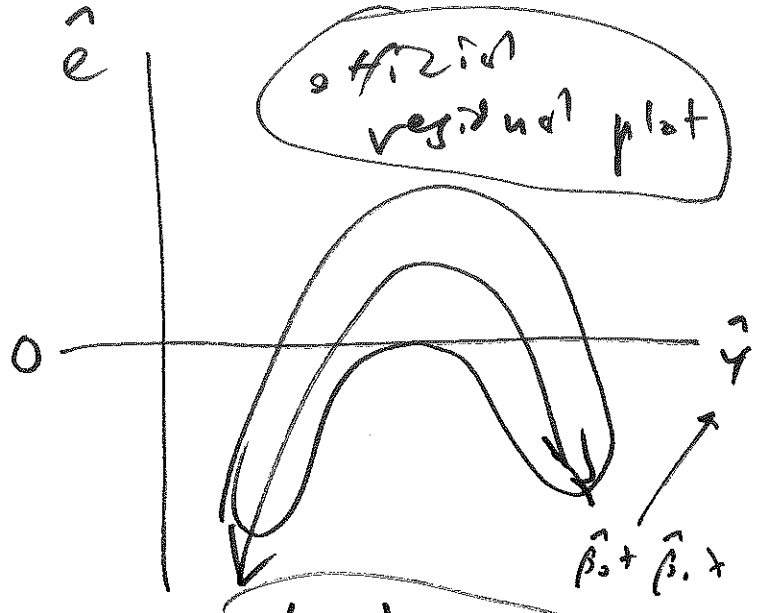
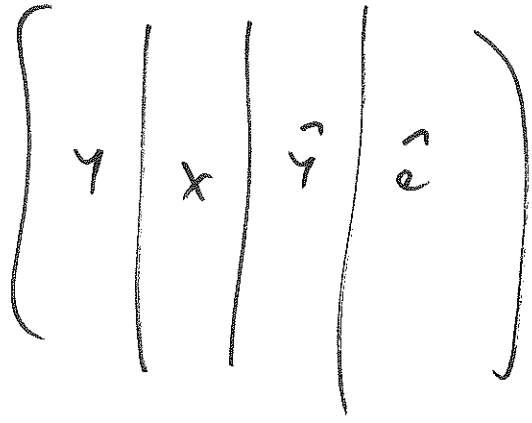
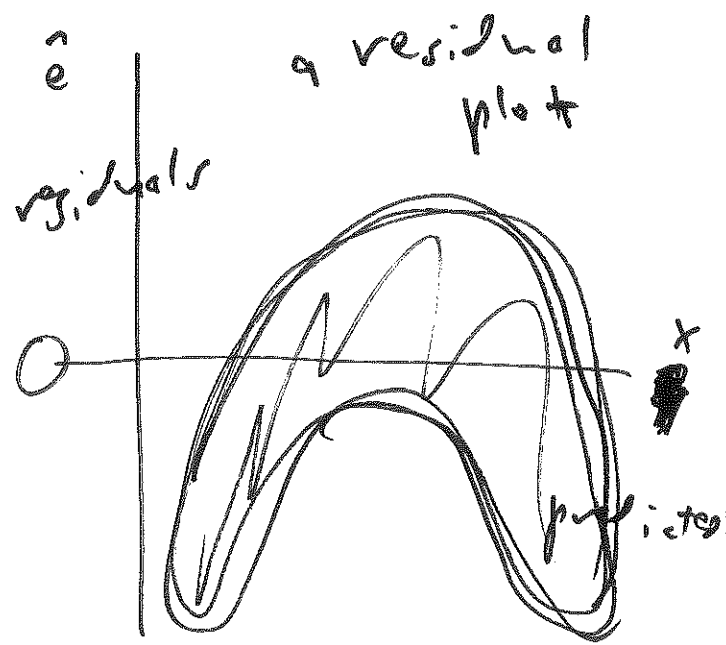
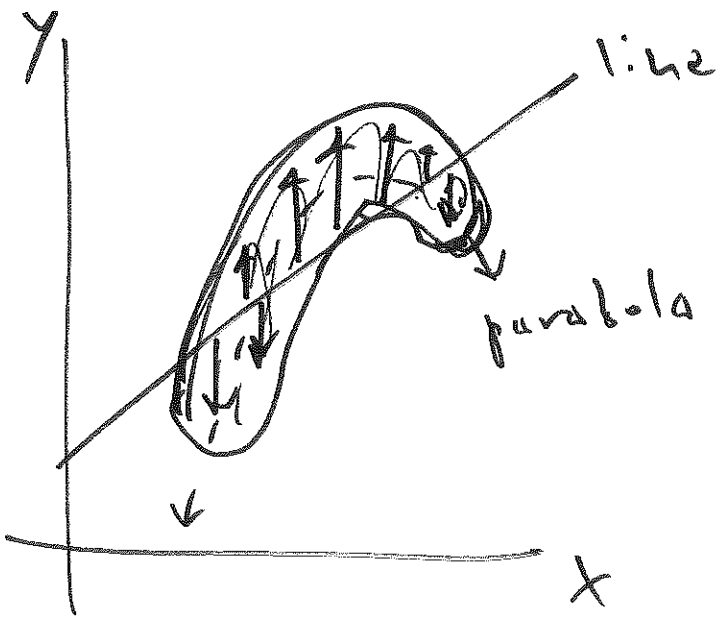
use x to predict y :

$$\hat{y}_{use\ x} = \hat{\beta}_0 + \hat{\beta}_1 x$$

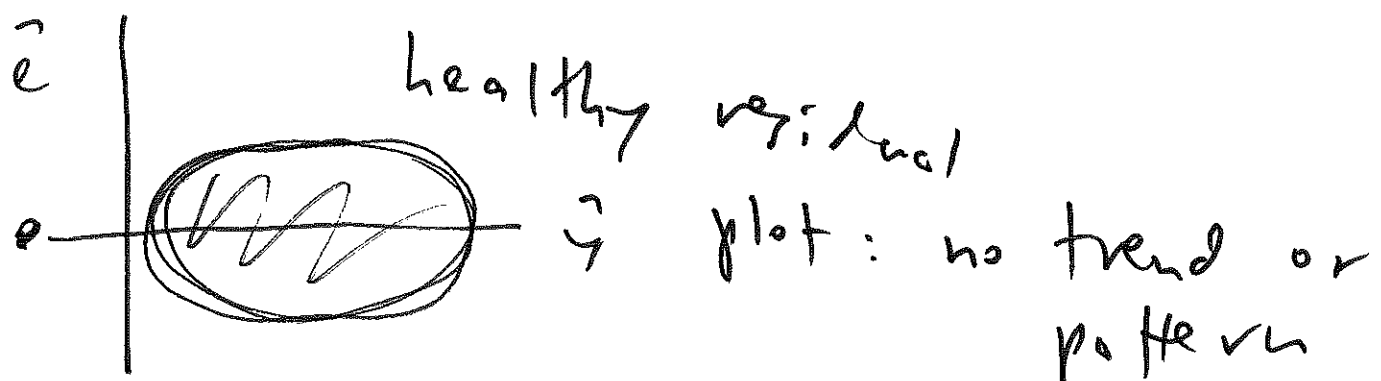
with $\hat{SE}(\hat{y}_{use\ x}) = s_y \sqrt{1 - r^2}$

this is smaller than $\hat{SE}(\hat{y}_{no\ x})$





unhealthy: nonlinearity (curvature)



$y = \text{height}$

$x = \text{trunk diameter}$

$y_i = \beta_0 + \beta_1 x_i + e_i$ (1)

Simple linear regression

$y = \text{height}$

$x_1 = \text{trunk diameter}$

\vdots

$x_k = x_k = \text{volume}$

multiple linear regression
 $k > 1$ variables

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + e_i$$

can generalize least squares to get estimates:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik}$$

$$\begin{array}{c}
 \left[\begin{array}{c|ccc}
 y_1 & x_{11} & x_{12} & \dots & x_{1k} \\
 y_2 & x_{21} & x_{22} & \dots & \vdots \\
 \vdots & \vdots & \vdots & \dots & \vdots \\
 y_n & x_{n1} & x_{n2} & \dots & x_{nk}
 \end{array} \right] \begin{array}{c} \vec{y}_1 \\ \vdots \\ \vec{y}_n \end{array}
 \end{array}$$

if
very good,
y & \vec{y}
should
be close

Q is the regression useful?
multiple

A could compute correlation between

y & $\vec{y} = R$; square to get

$R^2 =$ multiple $R^2 =$ coefficient of determination

↑ want big