

# Lecture notes to Stock and Watson chapter 4

## Introductory linear regression

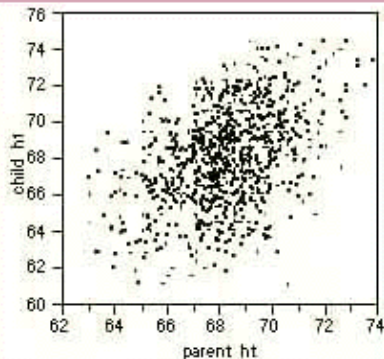
Tore Schweder

Sept 2009

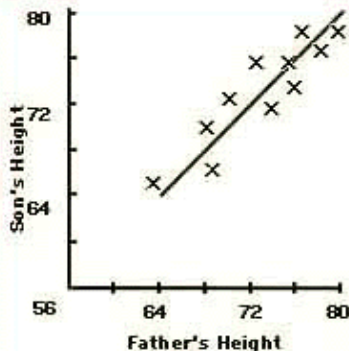
# Regression

"Regression" is due to Francis Galton (1822-1911): how is a son's height related to his father's height?

**Figure A**  
*Galton's Original Data*



*Hand-Drawn Data With  
Regression Line*



# Regression towards the mean

Plate X.

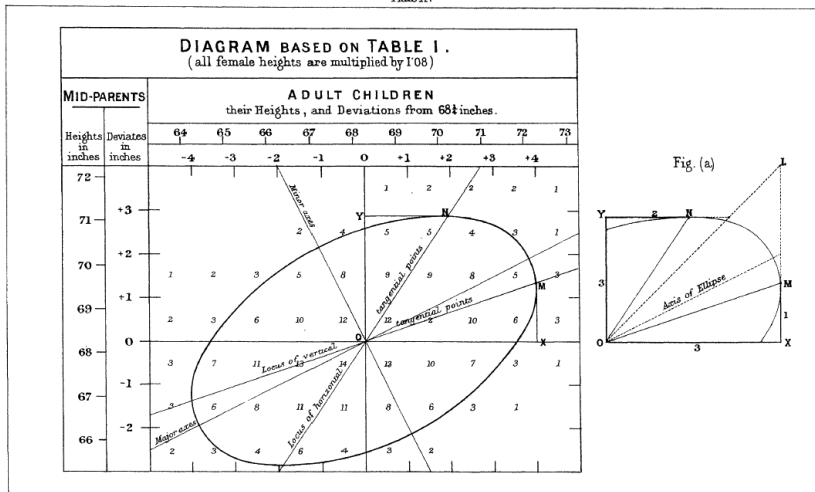
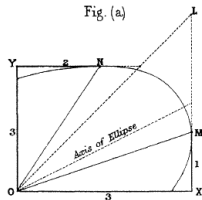


Fig. (a)



From Anthropological Inst., Vol. XV, Pl. X.

Figure: Galton's original diagram. Parents-child pairs by height.

# The geometry of linear regression

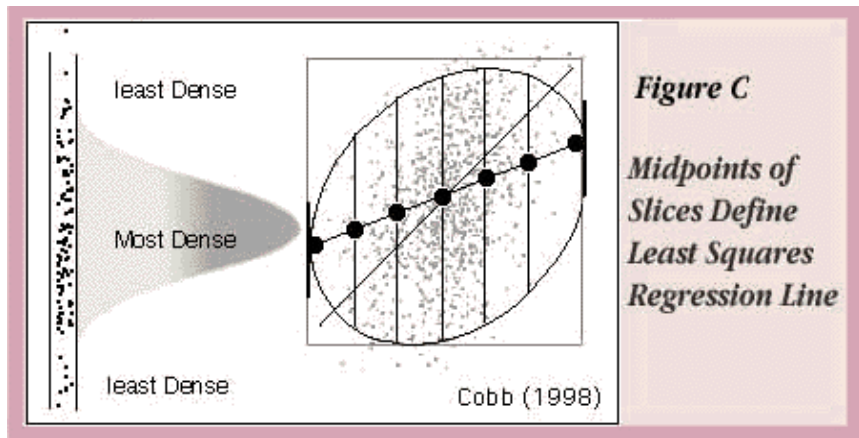


Figure: The regression curve is  $y = f(x) = E[Y|X = x]$ .

# The linear regression model

Question: how is a response variable  $Y$  related to a stimulus variable  $X$  (explanatory/control/explanatory)? Assuming a linear relation,

- 1 what is the slope?

# The linear regression model

Question: how is a response variable  $Y$  related to a stimulus variable  $X$  (explanatory/control/explanatory)? Assuming a linear relation,

- 1 what is the slope?
- 2 is the slope positive?

# The linear regression model

Question: how is a response variable  $Y$  related to a stimulus variable  $X$  (explanatory/control/explanatory)? Assuming a linear relation,

- 1 what is the slope?
- 2 is the slope positive?
- 3 how good does the estimated line fit the observed data?

# The linear regression model

Question: how is a response variable  $Y$  related to a stimulus variable  $X$  (explanatory/control/explanatory)? Assuming a linear relation,

- 1 what is the slope?
  - 2 is the slope positive?
  - 3 how good does the estimated line fit the observed data?
- Example:  $Y = \text{growth in BNP}$ ,  $X = \text{inflation previous year}$ .



# The linear regression model

Question: how is a response variable  $Y$  related to a stimulus variable  $X$  (explanatory/control/explanatory)? Assuming a linear relation,

- 1 what is the slope?
  - 2 is the slope positive?
  - 3 how good does the estimated line fit the observed data?
- Example:  $Y =$  growth in BNP,  $X =$  inflation previous year.
  - $D$ : A sample of size  $n$  of pairs explanatory variables  $X$  and response variable  $Y$ .

# The linear regression model

Question: how is a response variable  $Y$  related to a stimulus variable  $X$  (explanatory/control/explanatory)? Assuming a linear relation,

- 1 what is the slope?
  - 2 is the slope positive?
  - 3 how good does the estimated line fit the observed data?
- Example:  $Y =$  growth in BNP,  $X =$  inflation previous year.
  - $D$ : A sample of size  $n$  of pairs explanatory variables  $X$  and response variable  $Y$ .
  - $\mathbf{M}$ :  $(X_1, Y_1), \dots, (X_n, Y_n)$  is an iid random sample from an infinite population.  $Y = \beta_0 + \beta_1 X + u$ ;  
 $E[Y|X = x] = \beta_0 + \beta_1 x \Leftrightarrow Eu = 0, \quad cov(u, X) = 0. \beta_0, \beta_1$  are parameters.

Least squares and regression was known by Laplace and Gauss (80 years before Galton).

$$\begin{aligned}(\hat{\beta}_0, \hat{\beta}_1) &= \arg \min_{b_0, b_1} \sum_{i=1}^n (Y_i - (b_0 + b_1 X_i))^2 \\ \hat{\beta}_1 &= \frac{s_{XY}}{s_X^2} = r \frac{s_Y}{s_X}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}\end{aligned}$$

- Predicted response given the stimuli:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$  is the fitted value.

Least squares and regression was known by Laplace and Gauss (80 years before Galton).

$$\begin{aligned} (\hat{\beta}_0, \hat{\beta}_1) &= \arg \min_{b_0, b_1} \sum_{i=1}^n (Y_i - (b_0 + b_1 X_i))^2 \\ \hat{\beta}_1 &= \frac{s_{XY}}{s_X^2} = r \frac{s_Y}{s_X}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \end{aligned}$$

- Predicted response given the stimuli:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$  is the fitted value.
- (Empirical) residual:  $\hat{u}_i = Y_i - \hat{Y}_i$  is the vertical distance from the fitted value to the observed value of  $Y$ .  $\bar{\hat{u}} = 0$ ,  $r_{\hat{u}X} = 0$

- Homoscedasticity when  $\text{var}(Y|X = x) = \sigma^2$  is independent of  $x$ .

- Homoscedasticity when  $\text{var}(Y|X = x) = \sigma^2$  is independent of  $x$ .
- Then  $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n-2} \text{RSS}$  is unbiased.

- Homoscedasticity when  $\text{var}(Y|X = x) = \sigma^2$  is independent of  $x$ .
- Then  $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n-2} \text{RSS}$  is unbiased.
- $TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \text{RSS} + \text{ESS}$

- Homoscedasticity when  $\text{var}(Y|X = x) = \sigma^2$  is independent of  $x$ .
- Then  $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n-2} \text{RSS}$  is unbiased.
- $TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \text{RSS} + \text{ESS}$
- $R^2 = \frac{\text{ESS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}} = r_{XY}^2$



# Be aware of extreme observations!

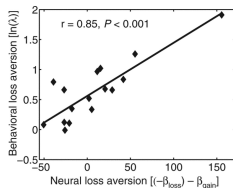


Figure: The OLS line hinges on one point.

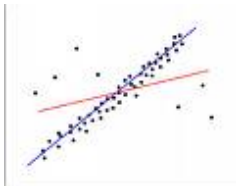
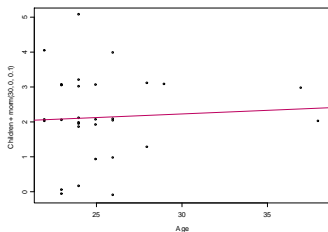


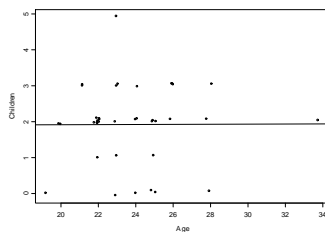
Figure: Heavy tails in the conditional distribution of  $Y|X$  distorts the OLS line.

# Desired number of children and age for beginning master students

	$n$	$\overline{Children}$	$\overline{Age}$	OLS	$\rho$	$R^2$
2008	38	1.92	23.8	$C = 1.88 + 0.002A$	0.004	0.00002
2009	29	2.13	25.4	$C = 1.61 + 0.021A$	0.06	0.004



2009



2008

# Desired number of children and number of siblings, 2009

```
. tabulate Children Siblings
```

Children	Siblings						Total
	0	1	2	3	4	6	
0	2	1	0	1	0	0	4
1	1	1	1	0	0	0	3
2	2	3	7	0	0	0	12
3	1	1	3	2	1	0	8
4	0	0	1	0	0	1	2
5	0	0	1	0	0	0	1
Total	6	6	13	3	1	1	30

```
. reg children siblings
```

source	SS	df	MS			
Model	8.93518887	1	8.93518887	Number of obs =	30	
Residual	34.5314778	28	1.23326706	F( 1, 28) =	7.25	
Total	43.4666667	29	1.49885057	Prob > F =	0.0110	
				R-squared =	0.2056	
				Adj R-squared =	0.1772	
				ROOT MSE =	1.1105	

children	coef.	std. err.	t	P> t	[95% conf. interval]	
siblings	.4214712	.156583	2.69	0.012	.1007255	.7422160
_cons	1.416832	.3346141	4.23	0.000	.7314064	2.102258

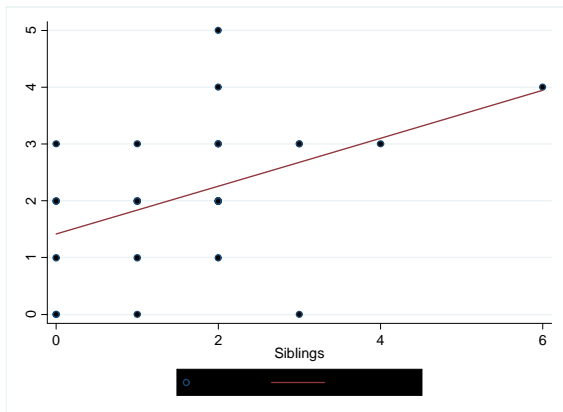
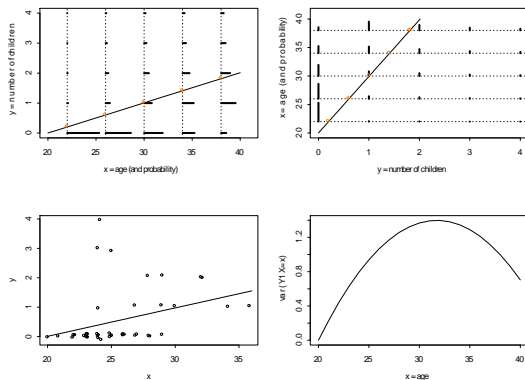


Figure: STATA scatter diagram of Children vs. Siblings, with OLS line. Some points represent several observations.

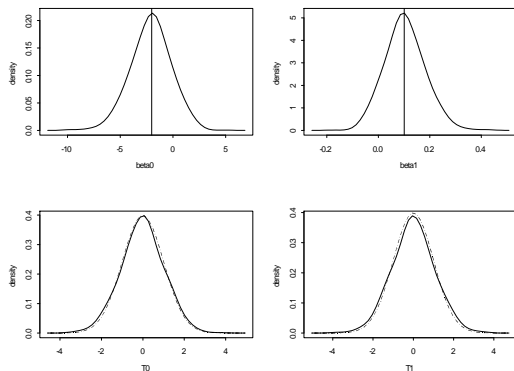
Is the (slightly) positive relation between Age and Children due to Siblings being an intermediate variable?  $\text{Age} \rightarrow \text{Siblings} \rightarrow \text{Children}$ , but controlling for siblings, no effect of Age on Siblings?

# Constructed model and simulated data



**Figure:** UL:  $E[Y = y|X = x] = -2.02 + 0.101x$ , and  $P[Y = y|X = x]$  shown by horizontal thick lines for some values of  $x$ ; UR: the same, but axes interchanged; LL the scatter plot of a simulated sample of size  $n = 43$ , the points are slightly jittered, with the OLS line  $y = -1.9 + 0.097x$ ; LR:  $\text{var}[Y|X = x]$ .

# Constructed model - repeated simulations



**Figure:** Approximate densities of  $\hat{\beta}_0$  (UL),  $\hat{\beta}_1$  (UR),  $T_0 = (\hat{\beta}_0 - \beta_0) / SE(\hat{\beta}_0)$  (LL), and  $T_1 = (\hat{\beta}_1 - \beta_1) / SE(\hat{\beta}_1)$  (LL), based on 1000 replicates of simulated data of size  $n = 43$  in constructed model for  $(X, Y)$ . Dotted curve is the normal density.

# Constructed model - repeated simulations (cont)

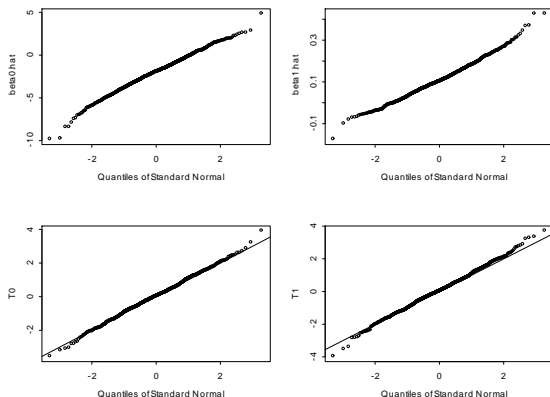


Figure: Normal probability plots (QQ-plots against  $N(0,1)$ ) for the same simulated material as in previous figure.

# Problems to be done in class

- SW: 4.3