

this correlation
time: & regression
next
time: analysis
of variance

read: LN pp.
L-214 + L-289

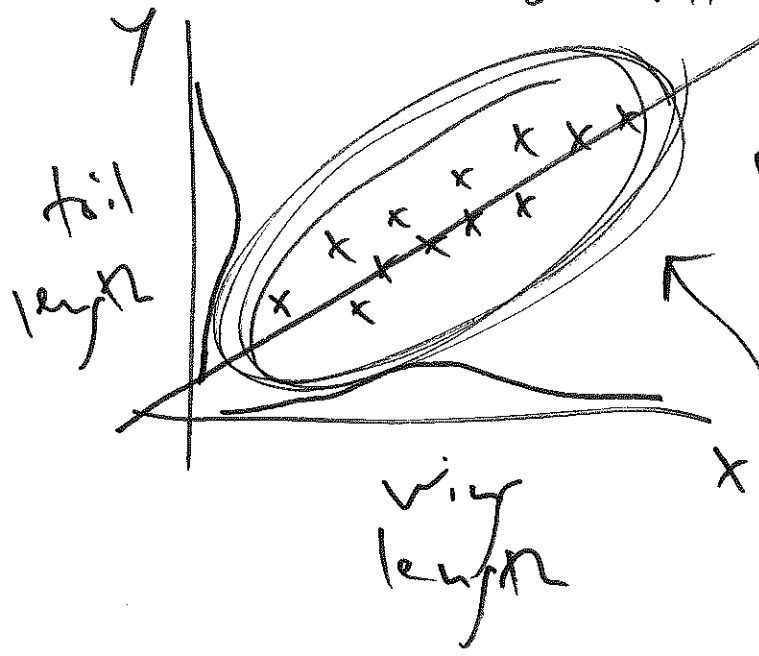
AMS7
15 Jul
2016

today: LN pp. L-214
→

please go to my uere & fill out the
online course evaluation by Fri 22 Jul

but 3 due Mon 18 Jul in class

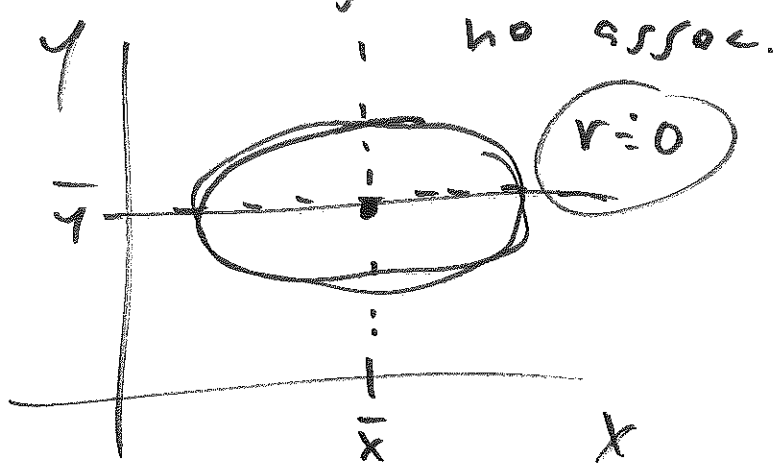
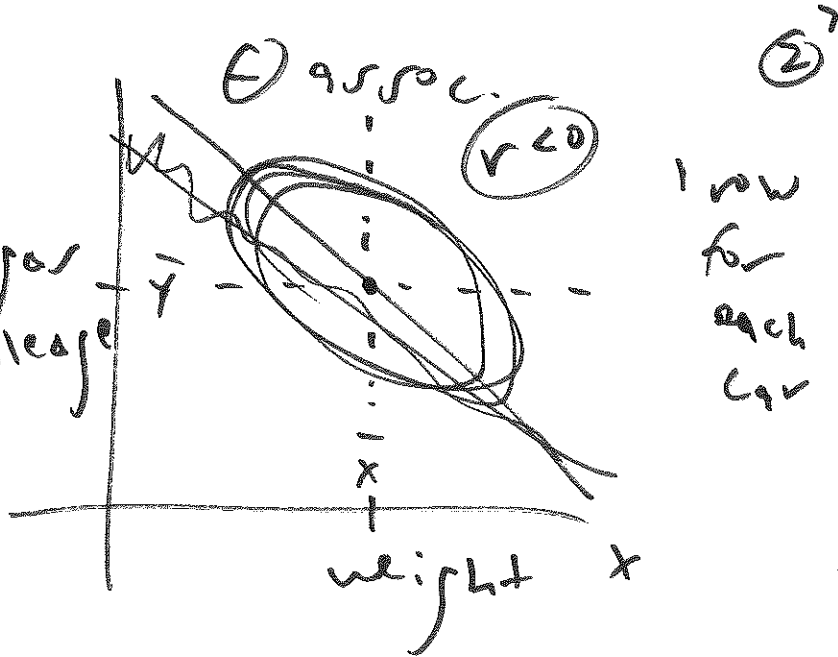
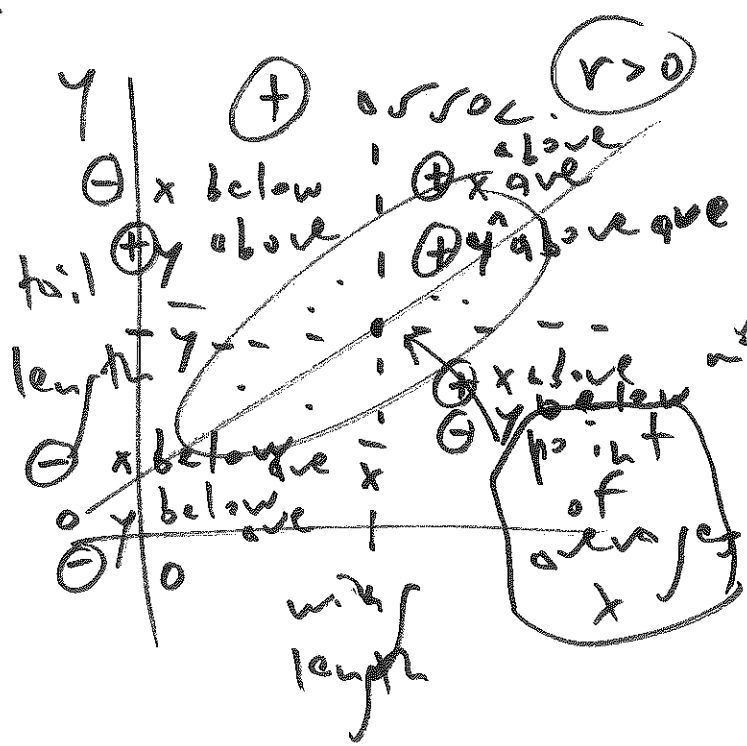
scatter plot (scatter
diagram)



elliptical
shape
(bivariate
normal dist.)
(x is normal,
y is normal)

Karl Pearson
(1890)

positive
association as $x \uparrow$,
 $y \uparrow$ on ave.



correlation (coefficient)
 $r =$ strength of linear association between x & y

y_1	x_1
y_2	x_2
\vdots	
y_n	x_n

$$r = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x^*} \right) \cdot \left(\frac{y_i - \bar{y}}{s_y^*} \right)$$

mean \bar{y} \bar{x}
 SD s_y s_x

$$s_x^* = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

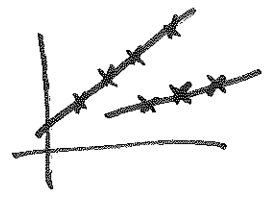
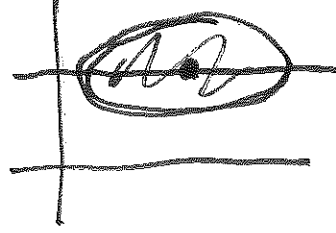
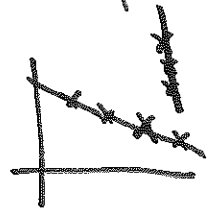
$$s_y^* = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

math facts about r

① r is a pure number (no units) ③

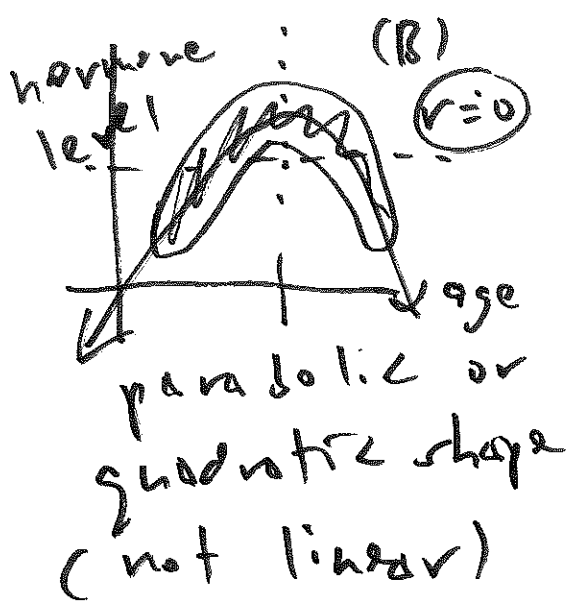
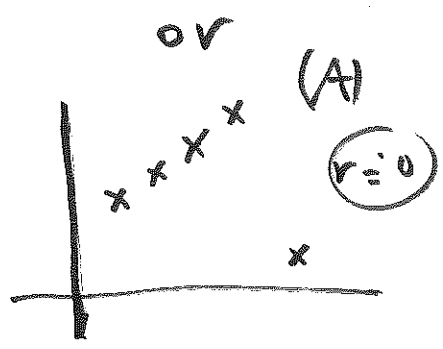


$-1 \leq r \leq +1$

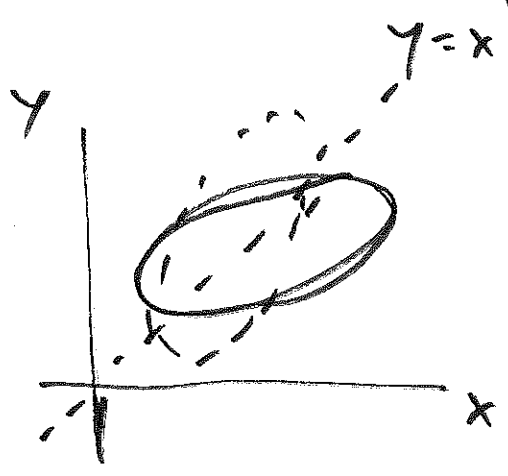


(9.55)

r can be fooled by (A) outliers or (B) non-linearity



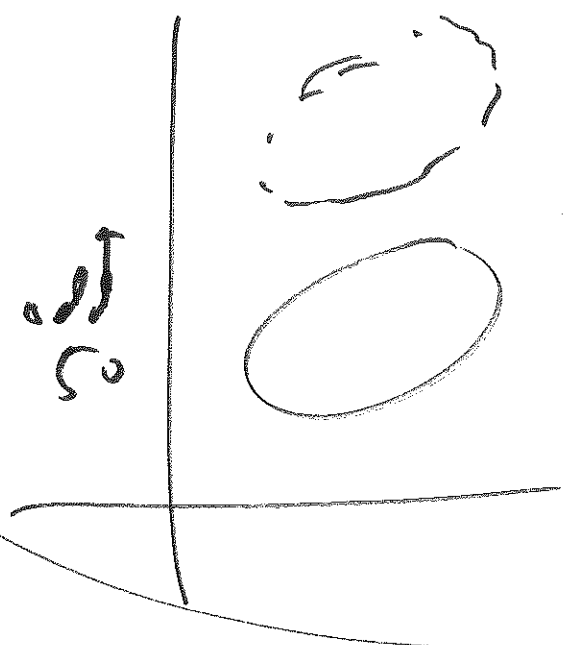
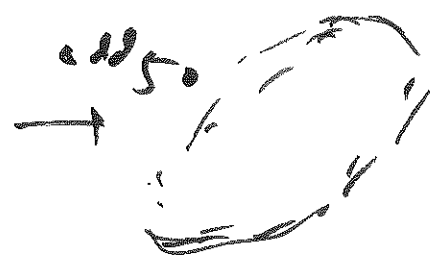
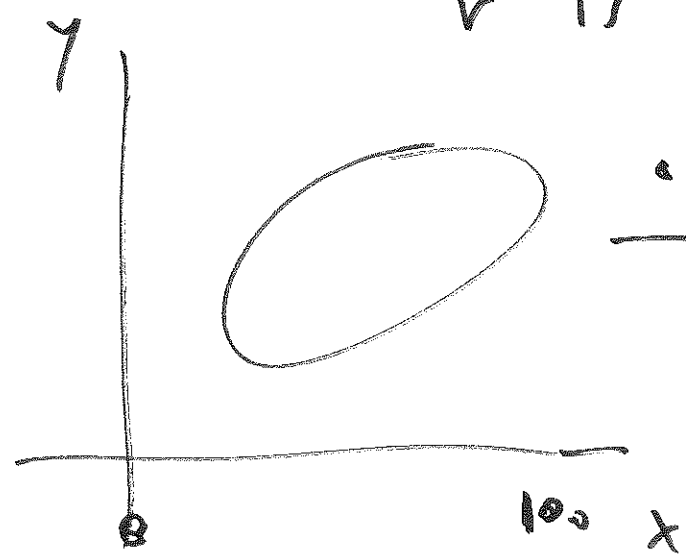
Spans:
 $r = +0.87$



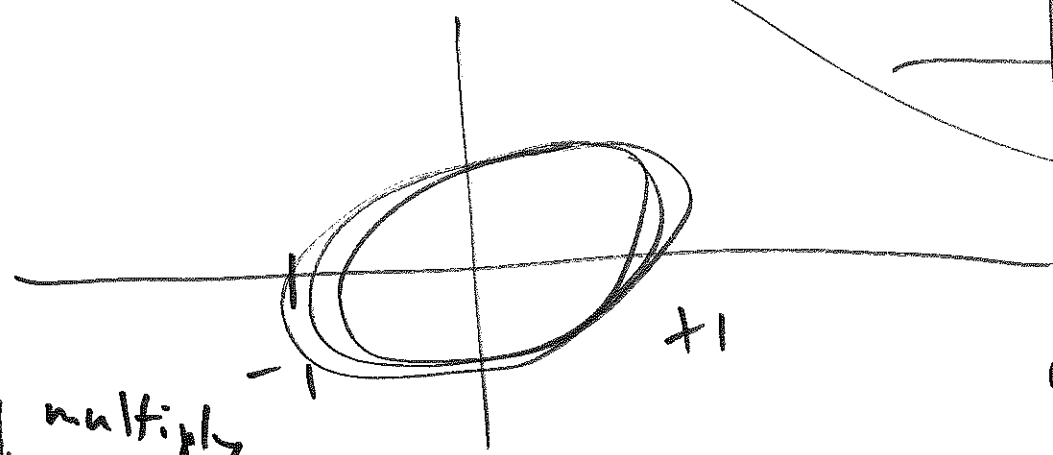
④ Corr. of any x with itself is +1

③ r is unchanged when the roles of x & y are reversed

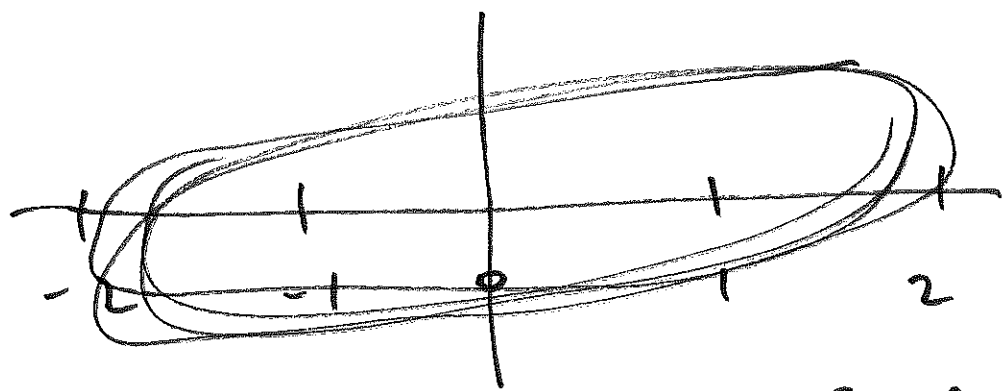
⑤ if you add (or subtract) a constant to (from) all the x values, r is unchanged; same with y



⑥



multiply by 2



if you multiply x by a positive constant, r stays the same; same with y

same with y

using r inferentially Q_1 is a corr. ⁽⁵⁾

between wing length (x) and tail length (y) of $r = +.87$ large in practical terms?

$A:$ smallest sparrows have

($x = 10 \text{ cm}$, $y = 7 \text{ cm}$); largest sparrows

have ($x = 11.5 \text{ cm}$, $y = 8.25 \text{ cm}$);

is the difference between 7 cm for

y & 8.25 cm large in practical

terms?

$$\frac{8.25 \text{ cm} - 7 \text{ cm}}{7 \text{ cm}} = \frac{1.25}{7} = 18\%$$

$Q_2:$ is this corr. of $+.87$

large in statistical terms?

null hypothesis value: 0 / model $\hat{\rho} = .87$

note facts

① $E_{IID}(r) = \rho$

② $SE_{IID}(r) = \sqrt{\frac{1 - \rho^2}{n - 2}}$

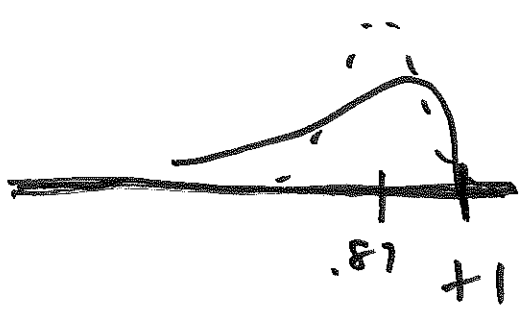
don't know

so use

$\hat{SE}_{IID}(r) = \sqrt{\frac{1 - r^2}{n - 2}}$

here $\hat{SE}(r) = \sqrt{\frac{1 - (+.87)^2}{12 - 2}} = 0.156$

≈ 0.16



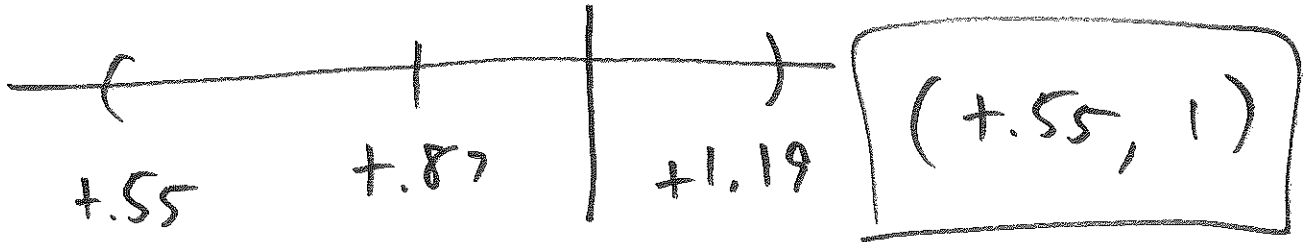
long run hist of r

let's try a large-sample (large n)

normal approximation: approx. 95% CI

for $\rho \rightarrow r \pm 1.96 \hat{SE}(r)$

here this becomes $(+.87) \pm 1.96 (.16)$ ⁷
 95% int. for ρ 0.32



t_1 truncate at t_1

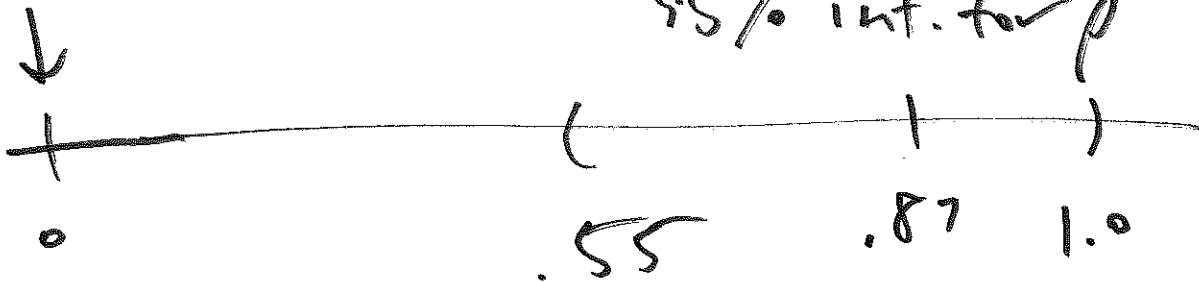
approx $(+.55, t_1)$

exact $(+.59, +.96)$

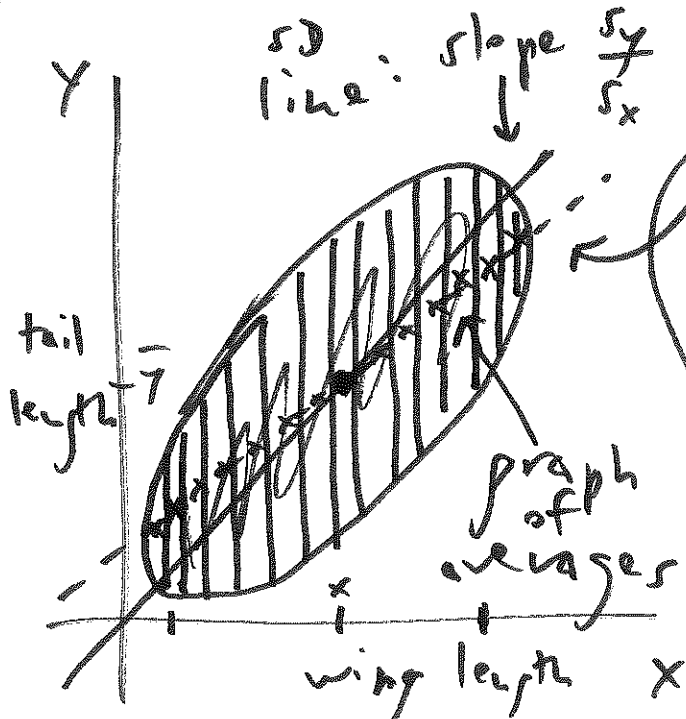
conservative
(wider)

LNL-231 - 244
 rr. optional

approx
95% int. for ρ



0 (way) not ⁱⁿ int; $r = +.87$ is highly
 statistically different from 0 (10.53)



regression line

Q: What's the equation of the line relating x & y ?

A₁: goal use line to predict y from x

A₂: goal use line to predict x from y

A₃: goal capture trend of $x-y$ relationship

A. equation of best line for predicting $y = \text{tail length}$ from $x = \text{wing length}$

Francis Galton (1890s)

here $\beta_1 = r \frac{s_y}{s_x} = \frac{+0.8704}{0.3950} \frac{0.3499 \text{ cm of tail length}}{\text{cm of wing length}} = 0.771$

slope of regression line

$$\hat{\beta}_1 = r \cdot \frac{s_y}{s_x}$$

math fact

⑨

(y)-intercept of reg. line

reg line has to go through (\bar{x}, \bar{y}) :

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

so

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

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

so

here y-int.

est. intercept
predicted y-value

est. slope

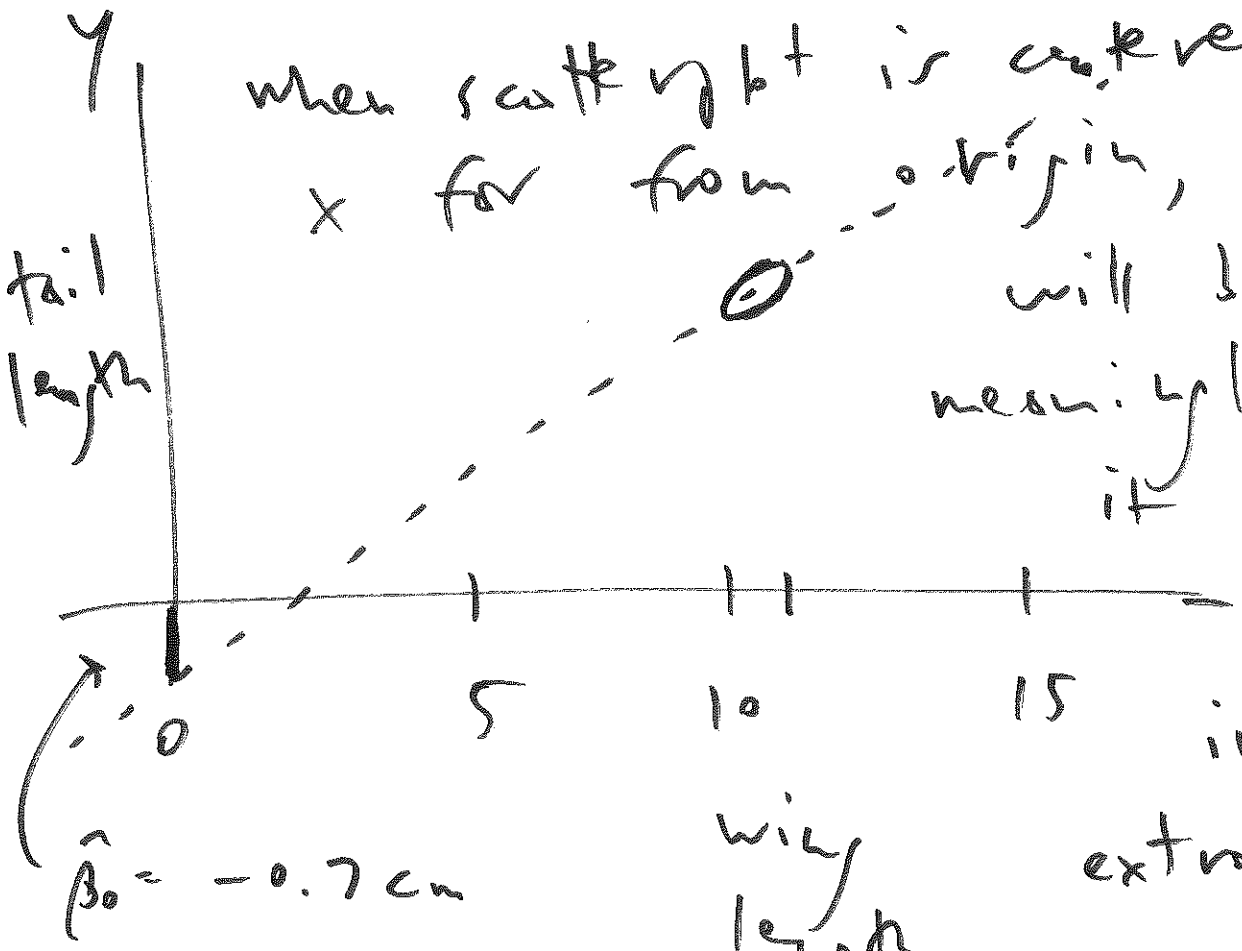
is

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = -0.669 \text{ cm of t.l.}$$

$$= 7.567 \text{ cm of tail length (+.1)} - \left(\frac{0.771 \text{ cm of t.l.}}{\text{cm of w.l.}} \right) (10.6 \text{ cm w.l.})$$

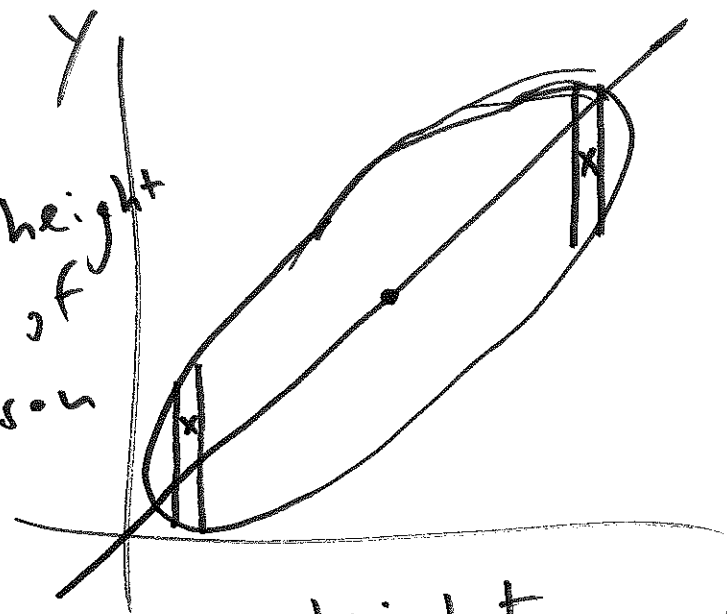
tail length

when scatter plot is centered in origin, y-intercept will be meaningless, because it involves



increase extrapolation away from the data

height of son



height of father

why regression?

tall (short) father tend to have

tall (short) sons, but not as tall (short) as the father