

Lecture 27 (ch. 11)

So far, we have made inferences about a model parameter, β . (α is in problem 11.10)

Q what about the true (pop.) prediction itself? $y(x) = \alpha + \beta x + \dots$

Unfortunately, the prediction $\hat{y}(x)$ has 2 different meanings:

- (point estimate of) the true/pop. conditional mean of y , given x . ← discussed last time
- (point) prediction of a single y , given x

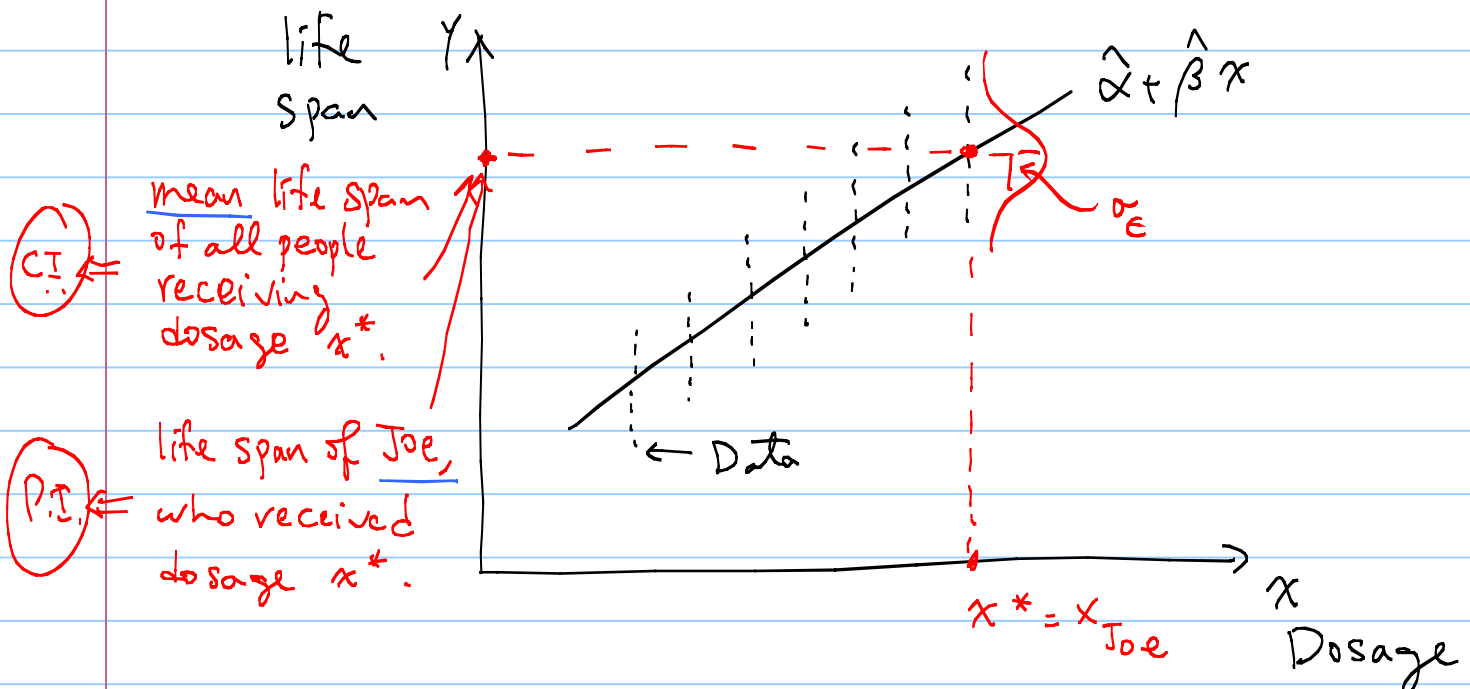
Note: The prediction $\hat{y}(x)$ is the same in both cases.

But the interpretation is different \Rightarrow different intervals & tests.

The two intervals/tests answer 2 diff. questions:

- \rightarrow What's the true cond'l mean of y for all cases, given $x = x^*$?
- \rightarrow What's the value of y for an individual case at $x = x^*$?

Example:



The first interval is just a confidence interval because it pertains to a pop. param (i.e. true mean of y , given x).

The 2nd interval is not a conf. interval at all!

It is called a Prediction Interval (P.I.).

The "levels" of the two intervals are often called
confidence level and prediction level

1) C.I. for The population mean, $y(x)$, given x :

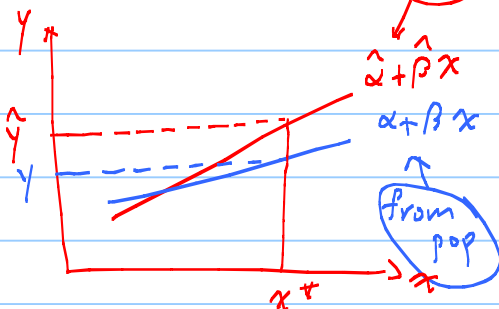
We need the sampling distr. of $\hat{y}(x)$.

The sampl. distr. of $\hat{y}(x) = \hat{\alpha} + \hat{\beta}x$ is Normal with params:

Thm. $\mu = y(x) = \alpha + \beta x$, $\sigma^2 = \sigma^2_{\text{estimation error}}$

where

estimation error = $\hat{y}(x) - y(x)$
 $\hat{y}(x) = \hat{\alpha} + \hat{\beta}x$ (sample fit)
 $y(x) = \alpha + \beta x$ (pop. fit)



$$\sigma^2_{\text{est. err}} = V[\text{est. err}] = V[\hat{y}(x)] + V[y(x)]$$

$$\sigma^2_{\text{est. err}} = \sigma^2_{\hat{y}} + 0 = 0$$

Approximate/estimate by

$$S^2_{\text{est. err.}} = S^2_{\hat{y}} + 0$$

$$\Rightarrow S^2_{\text{est. err.}} = S^2_{\hat{y}}$$

where $S^2_{\hat{y}} = S^2_e \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}} \right]$ No proof!

It follows That

$$z = \frac{\hat{y}(x) - y(x)}{\sigma_{\text{est. err}}} \sim N(0, 1), \quad \text{estimation error}$$

$$t = \frac{\hat{y}(x) - y(x)}{S_{\text{est. err}}} \sim t\text{-dist.} \quad df = n - 2 \quad \leftarrow k+1$$

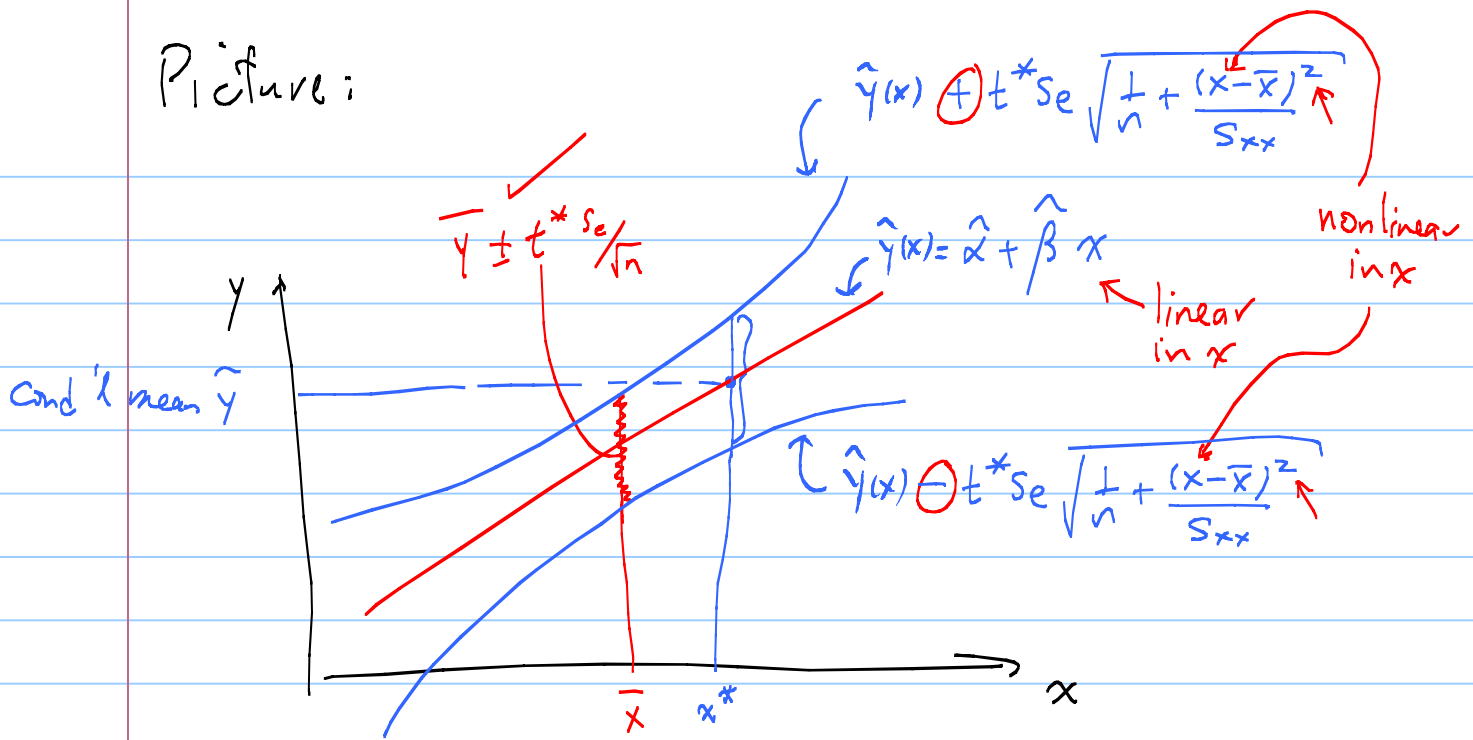
I.e. C.I. for mean $y(x)$, given x :

Table IV

$df = n - 2.$

$$\hat{y}(x) \pm t^* S_{\text{est. err}} = \hat{y}(x) \pm t^* S_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \quad \leftarrow k+1$$

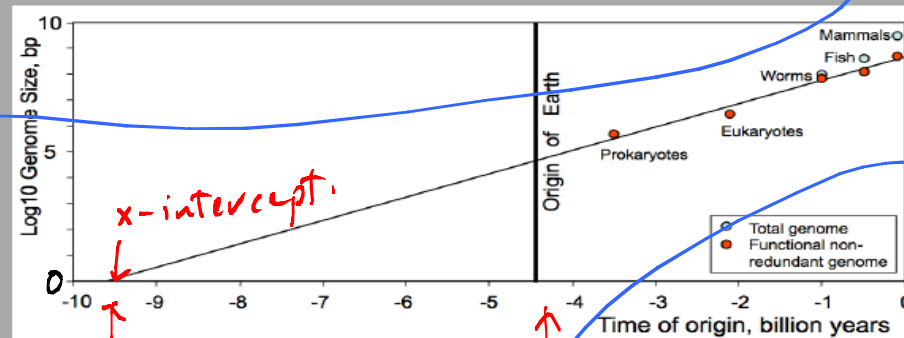
Picture:



Note: The C.I. gets wider the farther x gets from \bar{x} . Why? Regression has the property where the fit must go through the point $(x, y) = (\bar{x}, \bar{y})$. So, now, imagine a line that is fixed at that point. Any uncertainty in the slope will then cause the line to sweep a larger vertical direction in regions far away from $x = \bar{x}$.

The interpretation of these C.I.s is just as before.

What is most interesting in this relationship is that it can be extrapolated back to the origin of life. Genome complexity reaches zero, which corresponds to just one base pair, at time ca. 9.7 billion years ago (Fig. 1). A sensitivity analysis gives a range for the extrapolation of ± 2.5 billion years (Sharov, 2006). Because the age of Earth is only 4.5 billion years, life could not have originated on Earth even in the most favorable scenario (Fig. 2). Another complexity measure yielded an estimate for the origin of life date about 5 to 6 billion years ago, which is similarly not compatible with the origin of life on Earth (Jørgensen, 2007). Can we take these estimates as an approximate age of life in the universe? Answering this question is not easy because several other problems have to be addressed. First, why the increase of genome complexity follows an exponential law instead of fluctuating erratically? Second, is it reasonable to expect that biological evolution had started from something equivalent in complexity to one nucleotide? And third, if life is older than the Earth and the Solar System, then how can organisms survive interstellar or even intergalactic transfer? These problems as well as consequences of the exponential increase of genome complexity are discussed below.



origin of life

Earth's age

"Life before Earth", Sharov & Gordon (2013)

Q1: what is The correct conclusion after The CI is made?

Let μ_E = true age of Earth

μ_L = " " " Life

A) There is evidence that μ_E is smaller than μ_L .

B) There is evidence that μ_E is larger than μ_L .

C) There is evidence that μ_E and μ_L are comparable.

D) None of The above.

After making The CI for y (and accounting for errors in we cannot tell, because The range of possible x-intercepts includes The Earth's age.

"Earth Before Life", Marzban & Yurtsever (2014)

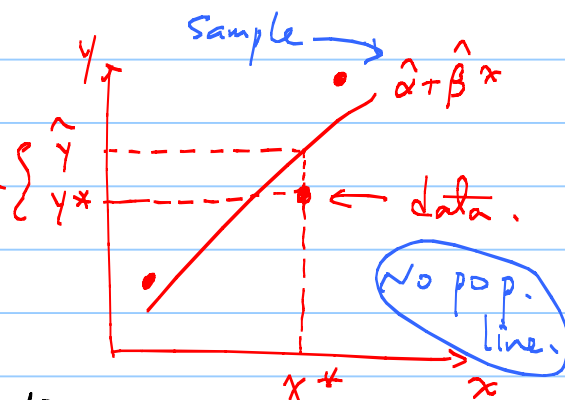
2) Prediction Interval (P.I.) for a single y .

Suppose y^* is Joe's y value corresponding to his x -value, x^* .

A theorem states that $(\hat{y}(x) - y^*)$ has a normal distr. with params

$$\mu = 0, \quad \sigma^2 = \sigma_{\text{prediction error}}^2$$

where prediction error $= \hat{y}(x) - y^*$



$$\sigma_{\text{pred. err}}^2 = V[\text{pred. err}] = \underbrace{V[\hat{y}(x)]}_{\sigma_{\hat{y}}^2} + \underbrace{V[y^*]}_{\sigma_{y^*}^2} = \sigma_{\hat{y}}^2 + \sigma_{y^*}^2$$

$$\sigma_{\text{pred. err}}^2 = \sigma_{\hat{y}}^2 + \sigma_{\epsilon}^2$$

$$\therefore S_{\text{pred. err}}^2 = \underbrace{S_{\hat{y}}^2}_{\text{above}} + \underbrace{S_{\epsilon}^2}_{\text{SSE}/(n-2)} = \frac{S_{\hat{y}}^2}{k+1} + \frac{S_{\epsilon}^2}{n-2}$$

Think! The variance of all y values at a given x is σ_{ϵ}^2 .

$$Z = \frac{\boxed{\hat{y}(x) - y^*}}{\sigma_{\text{pred. err}}} \sim N(0, 1)$$

prediction error

$$t = \frac{\hat{y}(x) - y^*}{S_{\text{pred. err}}} \sim t\text{-distr. df} = n-2$$

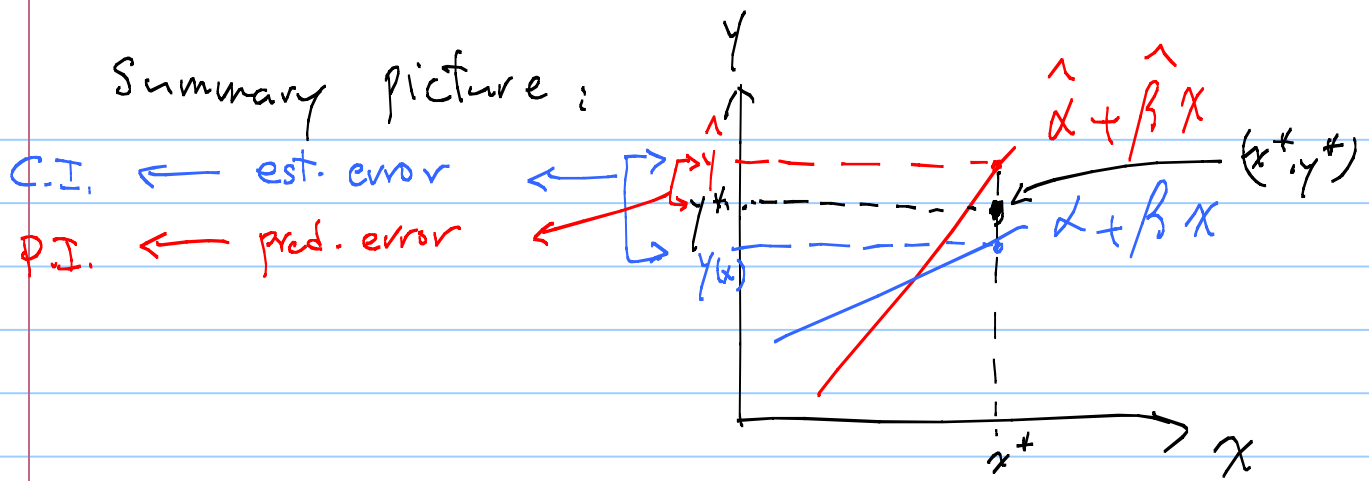
$k+1$

$$\therefore \text{P.I. for a single } y: \hat{y} \pm t^* S_{\text{pred. err}} = \hat{y} \pm t^* \sqrt{S_{\hat{y}}^2 + S_{\epsilon}^2}$$

Compare with C.I for y (The cond'l mean): $\hat{y} \pm t^* S_{\hat{y}}$

Q which one is bigger? P.I. makes sense?

Summary picture:



Don't forget what these intervals mean:

2 interpretations for C.I.:

- 1) We are 95% confident that the true (conditional) mean of y , given x , is in the observed C.I.
- 2) About 95% of random C.I.s will cover the true cond'l mean of y , given x .

For P.I. The most straightforward interpretation is

- 1) About 95% of random P.I.s will cover a single y , at a given x .

1') After we are more comfortable with interpretations we will allow ourselves to also say things like "plausible y values, at a given x , are in the observed P.I., at the 95% prediction level."

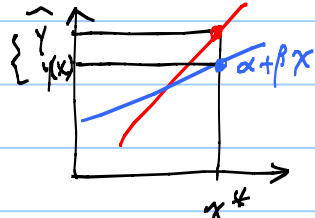
(See example, below)

CI, PI Side-by-side

In summary:

est. error

$$= \hat{y} - y(x)$$



$$\sigma_{\text{est. err}}^2 = \sigma_{\hat{y}}^2 + \sigma_{y(x)}^2$$

Recall that $\sigma_{y(x)}^2$ means the variance of $y(x)$ under resampling.

But $y(x)$ is the fit to the pop.

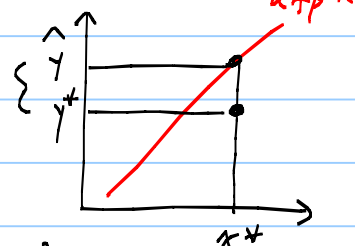
$$\text{So, } \sigma_{y(x)}^2 = 0.$$

$$\therefore \sigma_{\text{est. err.}}^2 = \sigma_{\hat{y}}^2$$

$$\therefore S_{\text{est. err.}}^2 = S_{\hat{y}}^2$$

pred. err.

$$= \hat{y} - y^*$$



$$\sigma_{\text{pred. err}}^2 = \sigma_{\hat{y}}^2 + \sigma_{y^*}^2$$

Again, $\sigma_{y^*}^2$ means the var. of y^* under resampling. But

y^* is the y for a given x , and so, its variance under resampling is just σ_e^2 .

$$\therefore \sigma_{\text{pred. err}}^2 = \sigma_{\hat{y}}^2 + \sigma_e^2$$

$$\therefore S_{\text{pred. err}}^2 = S_{\hat{y}}^2 + S_e^2$$

CI: $\hat{y} \pm t^* S_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$

$$S_{\hat{y}} \sim \frac{1}{\sqrt{n}} \rightarrow 0 \text{ as } n \rightarrow \infty$$

P.I.: $\hat{y} \pm t^* \sqrt{S_{\hat{y}}^2 + S_e^2}$

But PI does not!

One more comparison: How do CI & PI vary as $n \rightarrow \infty$?

I'll mention this tomorrow.

Example

11.20 (re-warded and revised, for clarity)

x = temperature

y = oxygen diffusivity.

$$n = 9, \sum x = 12.6 \quad \sum y = 27.68$$

$$\sum x^2 = 18.24 \quad \sum y^2 = 93.3448$$

$$\sum xy = 40.968$$

predict oxyg. diffusivity when temperature is 1.5 (in 1000 F°)
in a way that conveys info about reliability & precision.

11.5

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - n \bar{x}^2 = 18.24 - 9 \left(\frac{12.6}{9} \right)^2 = 0.6$$

$$SST = S_{yy} = \sum (y_i - \bar{y})^2 = \sum y_i^2 - n \bar{y}^2 = 93.3448 - 9 \left(\frac{27.68}{9} \right)^2 = 8.213$$

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - n \bar{x} \bar{y} = 40.968 - 9 \left(\frac{12.6}{9} \right) \left(\frac{27.68}{9} \right) = 2.216$$

$$\hat{y} = -2.095 + 3.6933 x$$

$$s_e = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{SST - \hat{\beta}(S_{xy})}{n-2}} = \sqrt{\frac{8.2134 - 3.6933(2.216)}{9-2}} = 0.0644$$

When temp = 1.5 in (1000 F°), what is the prediction
for the mean of diffusivity at that temp.?

A point estimate for that mean is given by the OLS line:

$$\hat{y} = \hat{\alpha} + \hat{\beta} x = -2.095 + 3.6933 x$$

$$\text{ie. } \hat{y} = -2.095 + 3.6933(1.5) = 3.445$$

A C.I. for the true mean at that temp. gives an interval estimate of that mean:

$$\hat{y} \pm t^* S_{\text{est. err}} = \hat{y} \pm t^* s_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$

$$= 3.445 \pm \underset{\substack{\uparrow \\ df = 9 - 2}}{2.365} (0.0644) \sqrt{\frac{1}{9} + \frac{(1.5 - \frac{12.6}{9})^2}{0.6}}$$

$$0.02302 = S_{\text{est. err}} = s_{\hat{y}}$$

\therefore C.I. for mean, $y(x)$, at temp = 1.5 : 3.445 ± 0.0544
(3.39, 3.50)

Interpretation?
See how, below.

for a single case

predict oxyg. diffusivity when temperature is $1.5\text{ K}^\circ\text{F}$
in a way that conveys info about reliability & precision.

this is asking for a prediction interval:

Interval estimate.

$$\begin{aligned} \hat{y} \pm t^* \sqrt{s_{\hat{y}}^2 + s_e^2} &= 3.445 \pm 2.365 \sqrt{(0.02302)^2 + (0.0644)^2} \\ &= 3.445 \pm 0.1617 = \underline{(3.28, 3.61)} \end{aligned}$$

- 1) 95% of such PI's will cover single observations of y at $x=1.5$
- 2) At 95% prediction level, plausible values for a single observation on y , at $x=1.5$, are between 3.28 and 3.61.

hw. lect 27-1

Give 2 interpretations (one involving confidence, the other involving probability) of the CI in above example.

hw. lect 27-12 Revised 11-18

Mist (airborne droplets or aerosols) is generated when metal-removing fluids are used in machining operations to cool and lubricate the tool and work-piece. Mist generation is a concern to OSHA, which has recently lowered substantially the workplace standard. The article "Variables Affecting Mist Generation from Metal Removal Fluids" (Lubrication Engr., 2002: 10-17) gave the accompanying data on x = fluid flow velocity for a 5% soluble oil (cm/sec) and y = the extent of mist droplets having diameters smaller than some value:

x :	89	177	189	354	362	442	965
y :	.40	.60	.48	.66	.61	.69	.99

a. Make a scatterplot of the data. By R.

b. What is the point estimate of the beta coefficient? (By R.) Interpret it.

c. What is s_e ? (By R.) Interpret it.

d. Estimate the true average change in mist associated with a 1 cm/sec increase in velocity, and do so in a way that conveys information about precision and reliability.

Hint: This question is asking for a CI for beta. Compute it AND interpret it.

By hand; i.e. you must use the basic formulas for the CI. E.g. for beta:

$\text{beta_hat} \pm t^* s_e / \sqrt{S_{xx}}$,

but you may use R to compute the various terms in the formula.

Use 95% confidence level.

e. Suppose the fluid velocity is 250 cm/sec. Find the mean of the corresponding y in a way that conveys information about precision and reliability. Use 95% confidence level. Interpret the resulting interval. By hand, as in part d.

f. Suppose the fluid velocity for a specific fluid is 250 cm/sec. Predict the y for that specific fluid in a way that conveys information about precision and reliability. Use 95% prediction level. Interpret the resulting interval. By hand, as in part d.

This document was created with Win2PDF available at <http://www.win2pdf.com>.
The unregistered version of Win2PDF is for evaluation or non-commercial use only.
This page will not be added after purchasing Win2PDF.