

what you have learned in this class.

Dealing with ambiguity

Random variable

histograms

Comparative boxplots

quantiles

distributions

probability (e.g. from Poisson)

sample mean and variance

distr mean and variance

qqplots

scatterplots

correlation

regression (multiple, polynomial, ...)

ANOVA (R^2 , $s_e \sim RMSE$)

overfitting, collinearity, interaction

sampling distribution

1-sample Confidence Interval for ...

2-sample CI for ...

t-distribution

Hypothesis testing with p-values

1-sample, 2-sample, paired, ... tests

tests for means and proportions

chi-squared test of multiple proportions in 1 pop

chi-squared test of independence of two categorical variables ← skipped.

1-way ANOVA F-test for the equality of multiple pop means.

t-test of regression coefficients

Confidence and Prediction Intervals

F-test of model utility

Model selection via bootstrapping (and cross-validation)

Neural networks (as a regression model).

) skipped.

Lecture 26 (Ch. 11)

problem 11.11

Last time, we did inference for 1 β (and α) t -interval & t -test
and for many β 's p -value " & F -test.

Q

what about the true (pop.) prediction itself? $y(x) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots$

Unfortunately, The (sample) prediction $\hat{y}(x)$ has 2 diff. interpretations.

-(point estimate of) The true/pop. conditional mean of y , given x , i.e. $y(x)$.

-(point) prediction of a single y , given x , call it y^* .

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

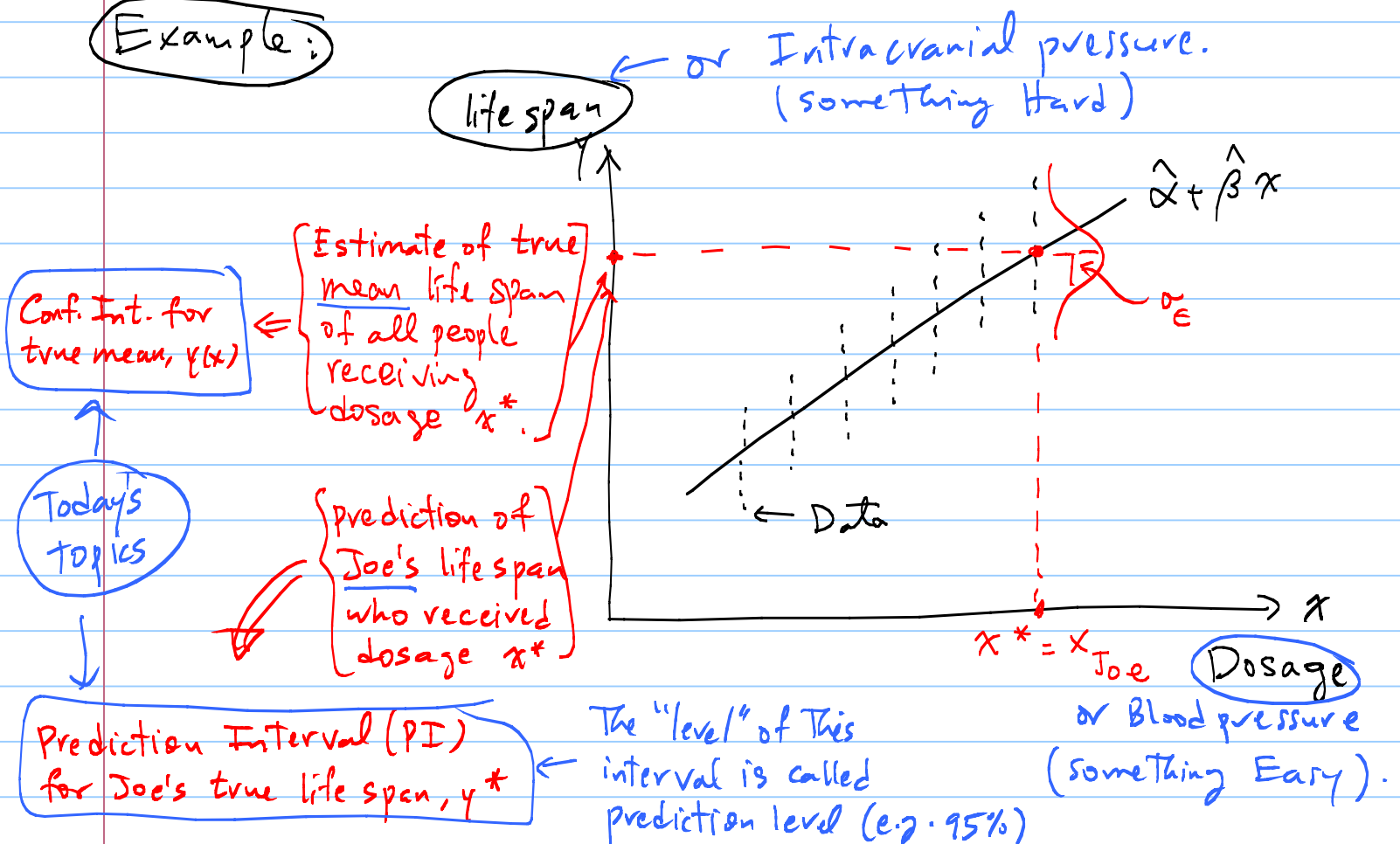
But the interpretation is different \Rightarrow different intervals & tests.

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 predicted y for an individual case at $x = x^*$?

Example:



(CI for $y(x)$)

- 1) For C.I. of The population mean, $y(x)$, given x : we need the sampling distr. of $\hat{y}(x)$.
- Analogous to:
For C.I. of μ_x , we need sampl. dist. of \bar{x} .
 $\bar{x} \sim N(\mu, \sigma^2/n)$

Theorem. The sampl. distr. of $\hat{y}(x) = \hat{\alpha} + \hat{\beta}x$ is Normal with params:
 $\mu = y(x) = \alpha + \beta x$, $\sigma^2 = \sigma^2_{\text{estimation error}}$

where

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

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

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

$$\sigma^2_{\text{est. err}} = \sigma^2_{\hat{y}} + 0 \quad \text{because } y(x) = \text{pop fit.}$$

Approximate/Estimate The σ 's with Their sample analog:

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

$$S^2_{\hat{y}} = S_e^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}} \right]$$

No proof

It follows That

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

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

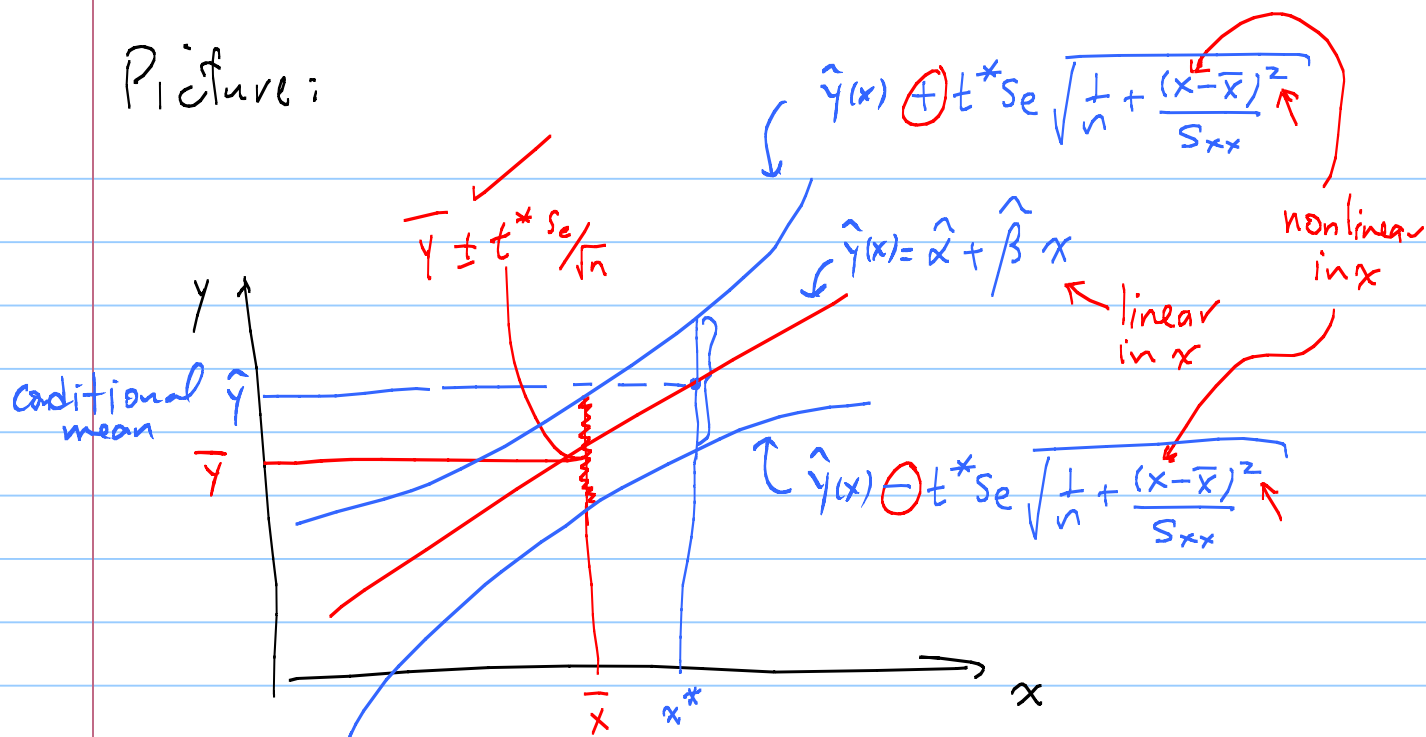
I.e. C.I. for mean $y(x)$, given x :

Table IV

$df = n - 2.$

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

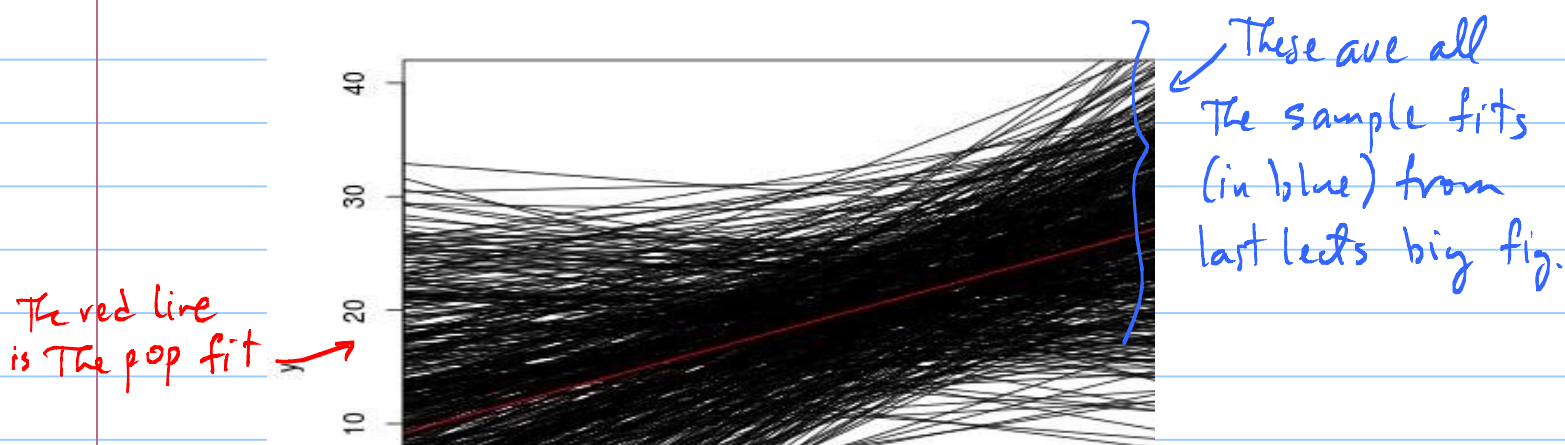
Picture:



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.

Recall, in "our" regression, \bar{x} = fixed (ie. x 's have no uncertainty). But \bar{y} is a r.v. (ie. every sample will have a different \bar{y}). So, that will shift the sample fit up/down across trials. BUT the $\hat{\beta}$ will also change across trials. So, now imagine what all the fits will look like:

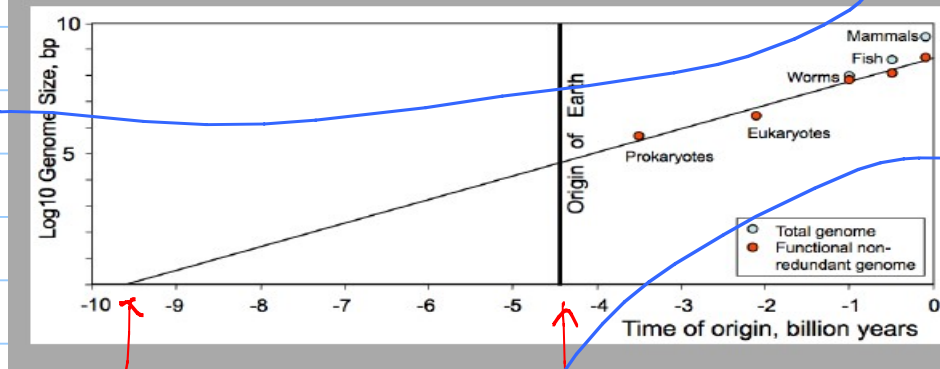


Each sample fit is forced to go thru (\bar{x}, \bar{y}) , and the slope changes. Do you see how a change in slope sweeps a wider region at larger x 's? That's the reason why the C.I. is wider for larger x 's.

C.I/P.I. are important, because they allow us to make Uncertainty "bands". Without them, wrong conclusions may follow. E.g.

Sharov & Gordon (2013) "Life Before Earth":

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.



Confidence Bands.

origin of life

Earth's age

From This

They conclude That Life predates Earth, and that life must have been formed on some other planet, then transported to Earth.

In a follow-up paper

(Marzban et al. (2014): "Earth Before Life", Biology Direct 9:1)

we showed that there are (at least) 2 problems with that analysis

1) Extrapolation is bad!

2) Confidence Intervals must be considered.

Don't forget what these intervals mean:

⇒ 2 interpretations for C.I.:

(Conditional) →

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

⇒ For P.I. The most straightforward interpretation is

- 1) About 95% of random P.I.s will cover y^* , at a given x .
- 2) We are 95% "confident" that y^* , at a given x , is in observed P.I.
(Technically, we should not use the word "confident" because that word is reserved for pop. params ($\mu, \pi, y(x), \dots$). So, people often say something like "plausible y values, at a given x , are in the observed P.I., at 95% prediction level."

(See example, below)

Note: The 3 "errors" in regression:

$$y_i = \alpha + \beta x_i + \epsilon_i$$

$$\left\{ \begin{array}{l} y - y(x) = \epsilon = \text{observation error} \Rightarrow \sigma_\epsilon \end{array} \right.$$

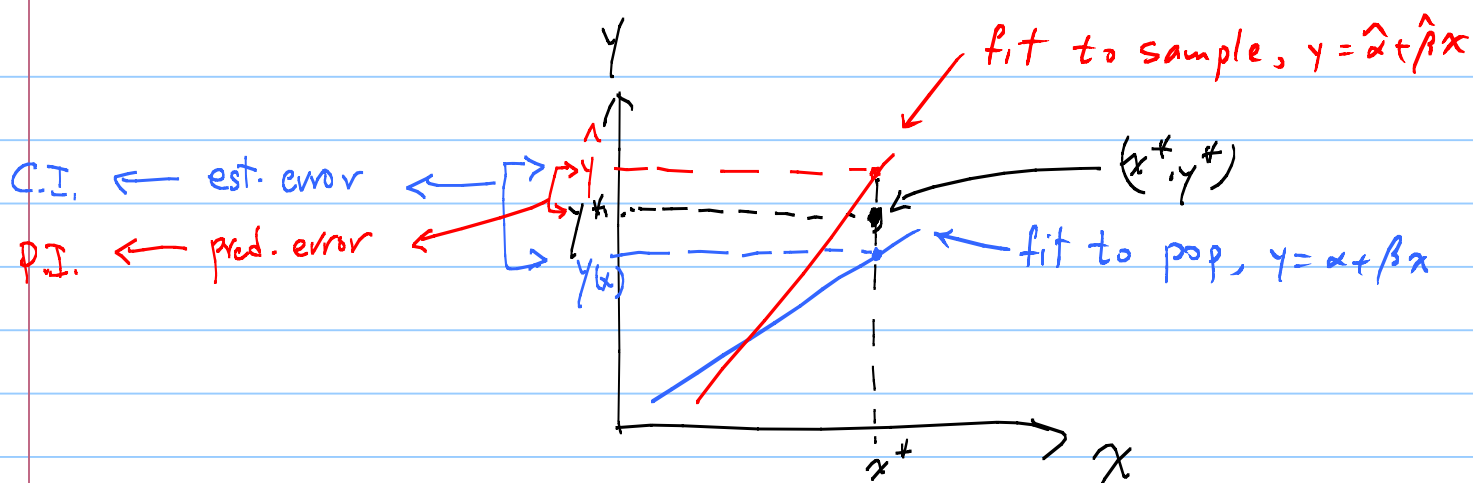
$$\left\{ \begin{array}{l} \hat{y}(x) - y(x) = \text{est. err} \Rightarrow \sigma_{\text{est. err}} \end{array} \right.$$

$$\left\{ \begin{array}{l} \hat{y}(x) - (y^*) = \text{pred. err.} \Rightarrow \sigma_{\text{pred. err}} \end{array} \right.$$

⌈ This is just a random y , at $x = x^*$.

You can even denote it as just y (without the $*$) as long as you remember that it's not observed.

CI, PI on top of each other



CI, PI side-by-side

est. error

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

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

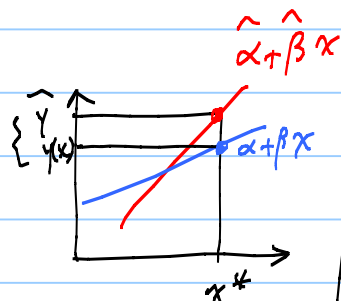
Recall that $\sigma_{y(x)}^2$ means the variance of $y(x)$ under resampling.

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. error

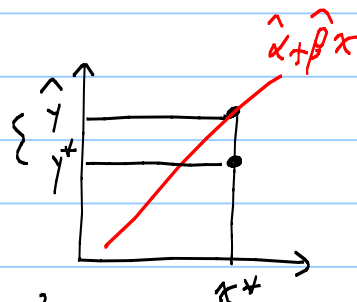
$$= \hat{y} - y^*$$

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

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 .

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

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



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

$S_{\hat{y}}$

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

Example

11.20 (re-warded and revised, for clarity)

x = temperature y = oxygen diffusivity.

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

$$\sum x^2 = 18.24 \quad \sum y^2 = 93.3448 \quad \sum xy = 40.968$$

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:

$$\begin{aligned}\hat{y} \pm t^* S_{\text{est. err}} &= \bar{y} \pm t^* s_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \\ &= 3.445 \pm \underset{\substack{\uparrow \\ df = 9-2}}{2.365} (0.0644) \sqrt{\frac{1}{9} + \frac{(1.5 - \frac{12.6}{9})^2}{0.6}} \\ &\quad S_{\text{est. err}} = s_{\hat{y}} = 0.02302\end{aligned}$$

\therefore (obs) CI for $y(x)$, i.e. mean of y at temp = 1.5 : 3.445 ± 0.0544
 $(3.39, 3.50)$

Interpretations: 1) With 95% confidence, the true mean of y ,
at $x = 1.5$, is between 3.39 and 3.50.
2) There is 95% prob that a random C-I will
cover the true mean of y at $x = 1.5$.

for a single case

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

this is asking for a prediction interval;

\uparrow
Interval
estimate.

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

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

hw-lect 26-1

In a simple regression problem, we have $n=16$, $\bar{x}=10$, $S_x = \frac{1}{\sqrt{8}}$, and $S_e = 4$. At $x=11$, what is the value of T_0 such that $\text{prob}(\text{prediction error} > T_0) = 0.01$?

(Hint: how do we standardize prediction error?)

hw-lect 24-2

Here is the regression version of a problem we have seen many a) times before. Show that, at a given x , the prob that a random y would fall into the obs CI for $y(x)$ is

$$\text{pr}\left(t_{obs} - t^* \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}}} < t < t_{obs} + t^* \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}}}\right)$$

$$\text{where } t_{obs} = \frac{\hat{y}_{obs}(x) - y(x)}{S_e}$$

(Hint: How do you standardize observation error?)

b) Show that, at a given x , the prob that a random $\hat{y}(x)$ would fall into the obs CI for $y(x)$ is $\text{pr}(t_{obs} - t^* < t < t_{obs} + t^*)$

$$\text{where } t_{obs} = \frac{\hat{y}_{obs}(x) - y(x)}{S_{\text{est. err.}}}$$

(Hint: how do you standardize estimation error?)

hw-optional

Consider The defining formulas for C.I and P.I :

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

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

$$\text{where } \hat{y}(x) = \hat{\alpha} + \hat{\beta}x$$

- a) As n becomes large (but not quite ∞) what does each of the following approach? For example $\hat{\alpha} \rightarrow \alpha$.

$\hat{\alpha} \rightarrow \alpha$ As n increases, $\hat{\alpha}$ approaches the population y -intercept α .

$\hat{\beta} \rightarrow$

$\hat{y}(x) \rightarrow$

$t^* \rightarrow$

$s_e \rightarrow$

$\bar{x} \rightarrow$

$S_{xx} = (n-1)S_x^2 \rightarrow$

- b) As $n \rightarrow \infty$, what does CI converge to?

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

- c) As $n \rightarrow \infty$, what does PI converge to?

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

hw_optional (revised 11.21)

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

- Make a scatterplot of the data. By R.
- What is the point estimate of the beta coefficient? (By R.) Interpret it.
- What is s_e ? (By R) Interpret it.
- 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.
- Suppose the fluid velocity is 250 cm/sec. Compute an interval estimate of the corresponding mean y value. Use 95% confidence level. Interpret the resulting interval. By hand, as in part d.
- 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.