

Index of 11-L

Page	Title
1	Practical information
2	Tips for comparing and presenting groups
3	Scatterplots and data example
4	Linear regression: Data + problem
5	Linear relation
6	Supplementary exercises 2.31 and 2.32
7	Linear regression model
8	Least squares estimation
9	Parameter estimates
10	Tests and confidence intervals
11	ANOVA Table for Linear Regression
12	Prediction
13	Standard errors in linear regression
14	How to report statistics in scientific papers?
15	Review: Another look at t and t^*

PRACTICAL INFORMATION

Schedule:

- home assignment III due today,
- home assignment IV on the website sometime next week
— absolutely final call for analysis of own dataset. . . .
- lab review this afternoon (note: 12:30-1:30 pm).

Today's lecture:

- ANOVAs revisited: brief review with two extra ideas,
- main topic of next two lectures: correlation and regression,
 - * today regression, including prediction,¹
 - * next week correlation + extra regression,²
 - * postpone references to r^2 , and regression vs. correlation until next week,
 - * entirely skip extra topics (IPS: scatterplot smoothers, nonlinear regression),
- additional topics:
 - * guidelines on how to report statistics in papers/theses,
 - * elaboration on t -values,
 - * links for choice of method added to media page.

¹ PSLS 3e: Chapters 4, 23; S: Sections 10.1-2; IPS 7e: Sections 2.1, 2.3, 10.1-2.

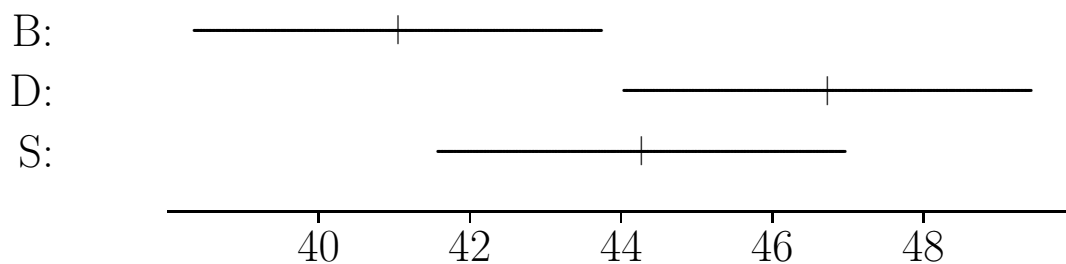
² PSLS 3e: Chapters 3, 23; S: Section 10.1; IPS 7e: Sections 2.2, 2.4, 10.2.

TIPS FOR COMPARING AND PRESENTING GROUPS

Group comparisons based on confidence intervals:

- based on test/CI for difference between parameters (e.g., $\mu_1 - \mu_2$),
- but conclusions available from group CIs in 2/3 cases (see figure):

Reading data example: 95% CIs for post3:



- * B vs. D: disjoint (non-overlapping) CIs \Rightarrow signif. ($P < 0.05$),
- * D vs. S: estimate in another CI \Rightarrow no signif. ($P > 0.05$),
- * B vs. S: need CI for difference ($\mu_B - \mu_S$) to assess signif.

assumes independent estimates, unadjusted for multiple testing.

Significance letter coding³: from software or constructed manually,

- order group means from lowest to highest,
- designate letter a to highest group + all groups not significantly different from it,
- designate letter b to next group in the same way (but drop if same pattern as for a),
- continue through all groups,
- Reading data: (uncorrected 5% error) coding: $B^b S^{ab} D^a$,
- Ex. 12.55 (vit A): (Bonferroni corrected) coding: $7^b 3^{ab} 1^{ab} 5^a 0^a$.

³ Meaning of letter codes: groups with same letter are *not* significantly different.

SCATTERPLOTS AND DATA EXAMPLE

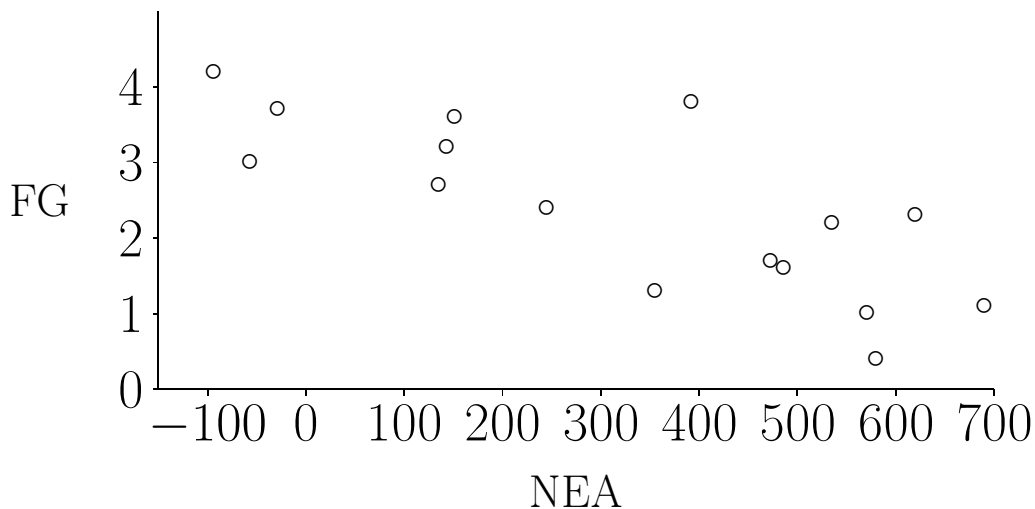
Scatterplot: a plot of two variables against each other:

- explanatory variable (if any) goes on the horizontal axis,
- one point per observation pair (Y, X) .

Data example: non-exercise activity (NEA) and fat gain (FG) in humans,⁴

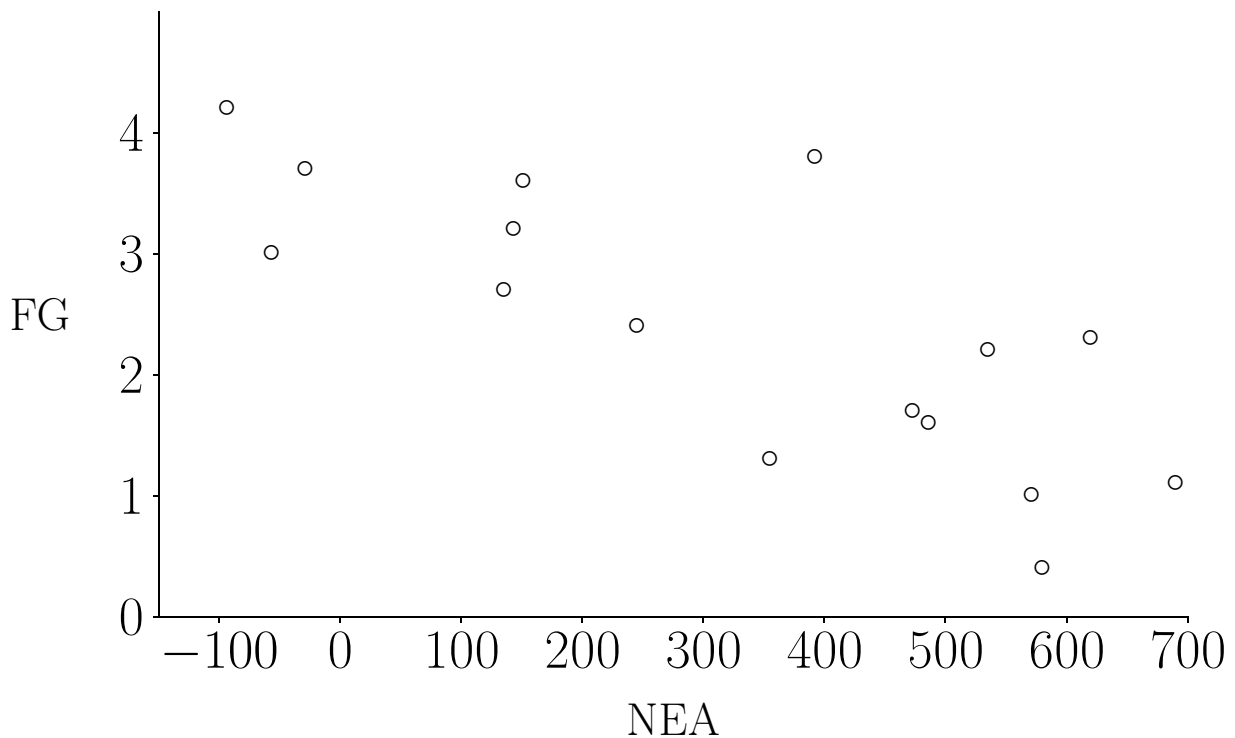
- for 16 young adults that were overfed for 8 weeks, measures of
 - * increase in NEA⁵, measured in calories,
 - * fat gain (FG), measured in kilograms,
- interest is in predicting fat gain from NEA,

- | | | | | | | | |
|----------|-----|-----|-----|-----|-----|-----|-----|
| fat gain | Y | 4.2 | 3.0 | 3.7 | 2.7 | ... | 1.1 |
| NEA | X | -94 | -57 | -29 | 135 | ... | 690 |



⁴ IPS 7e Example 2.18, data from Levine et al. (1999), *Science* **283**, 212-214.
⁵ NEA = any activity other than deliberate exercise, such as fidgeting, daily living, etc.; fidget(v): to make continuous small movements that annoy other people.

LINEAR REGRESSION: DATA + PROBLEM



Data:

$$\left. \begin{array}{l} Y_i = \text{fat gain} \\ X_i = \text{NEA} \end{array} \right\} \text{ for subject } i, i = 1, \dots, 16 = n.$$

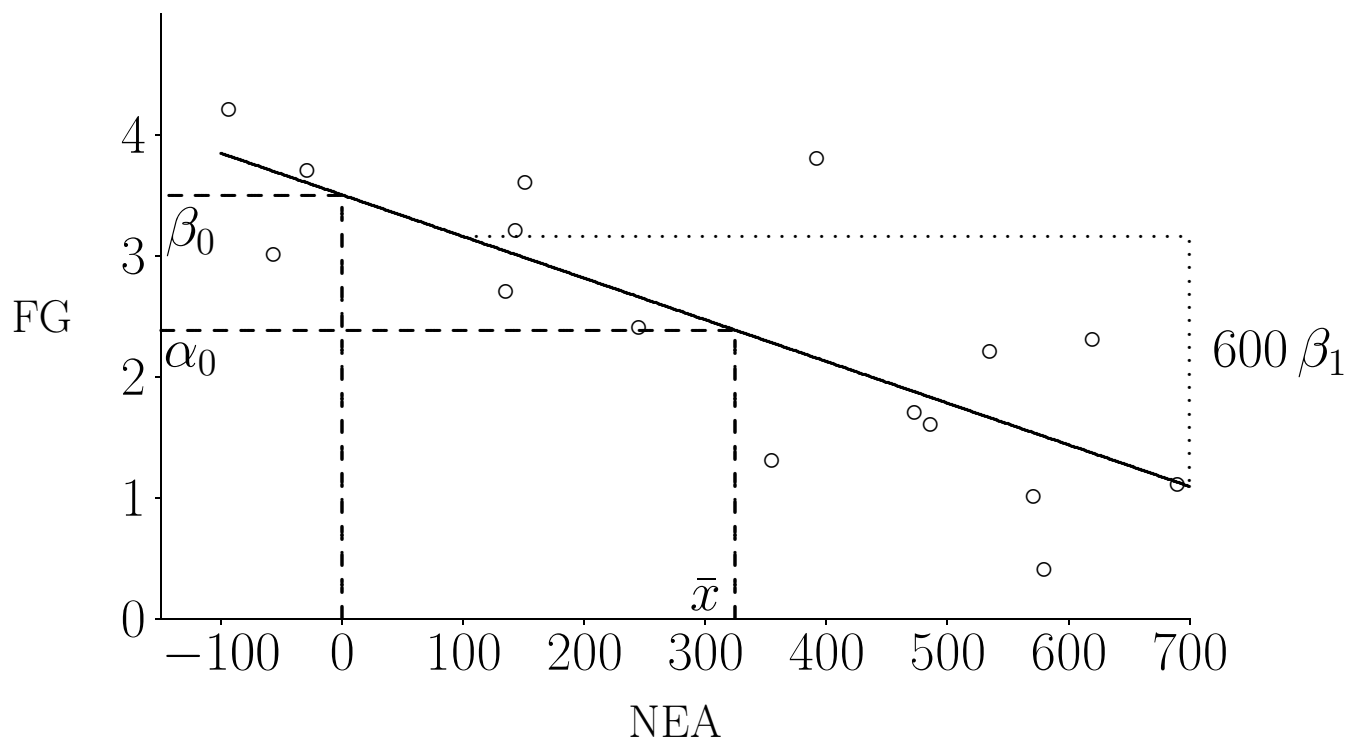
Problem: seek description of relationship between Y and X , in particular as: $y = f(x)$,

Why y as a function of x ?⁶

- causal relation? (if x controllable, we hope to impact y),
- interest in predicting y from x ?
(for prediction, x 's would be taken as fixed),
- X is not a random variable (\Rightarrow explanatory).

⁶ Commonly used (but somewhat imprecise) terminology to reflect this: y = dependent variable, x = independent variable.

LINEAR RELATION



Linear relation: $y = \beta_0 + \beta_1 \cdot x$

(or $y = a + bx$, as in Chapter 5 of PSLS, Chapter 2 of IPS):

- β_1 (or b) = slope of the line,
- β_0 (or a) = intercept of line with vertical axis ($x=0$),
- interpretation of slope: one unit increase in x implies a β_1 units change (increase or decrease) in y .

Alternative writing of the same line: $y = \alpha_0 + \beta_1(x - \bar{x})$:

- $\alpha_0 = y$ -value corresponding to $x = \bar{x}$,
- “centering” of x to avoid parameter (β_0) out of x ’s range,
- $\beta_0 = \alpha_0 - \beta_1 \bar{x}$, or $\alpha_0 = \beta_0 + \beta_1 \bar{x}$.

SUPPLEMENTARY EXERCISES 2.31 AND 2.32

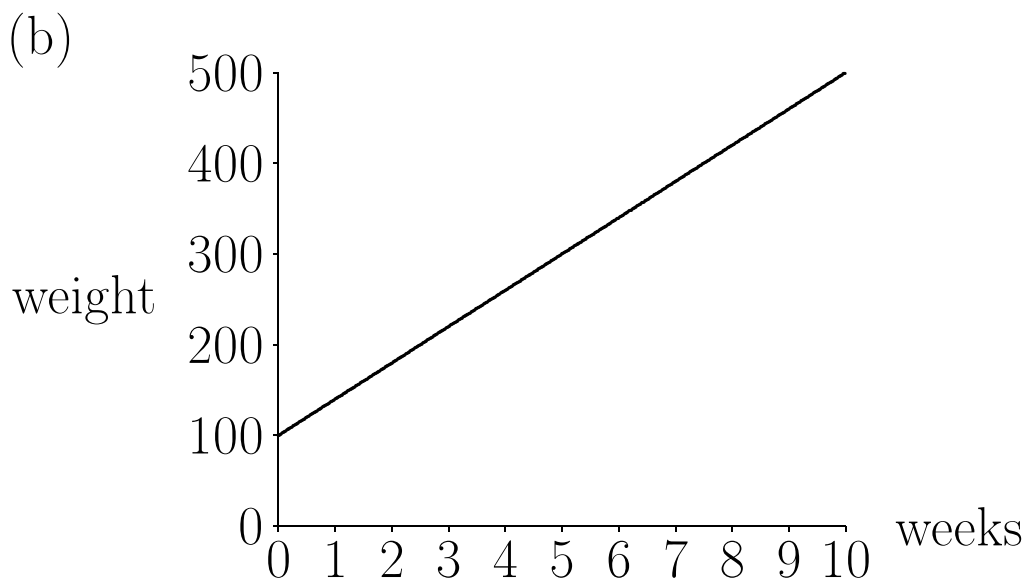
Exercise 2.31:

If x = number of seconds since splash and y = distance in meters, the equation is

$$y = 1500 (m/s) \cdot x (s).$$

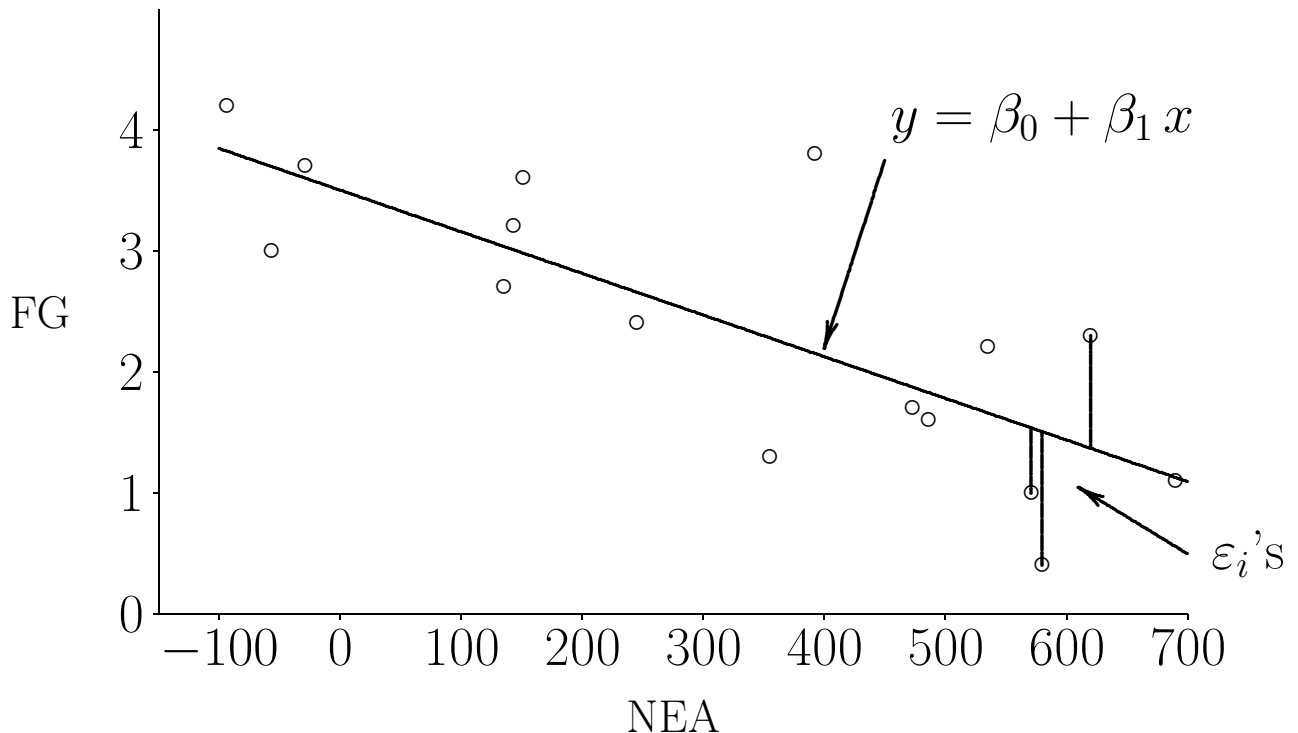
Exercise 2.32:

(a) equation: $\text{weight} = 100 + 40 \cdot \text{weeks}$,
and the slope of the line is $40 (g/\text{week})$.



(c) to use the linear equation for 2 years (104 weeks) would be an extreme case of *extrapolation* and clearly invalid, because rats do not continue to grow at a linear rate. Predicted value = $100 + 40 \cdot 104 = 4260 g$.

LINEAR REGRESSION MODEL



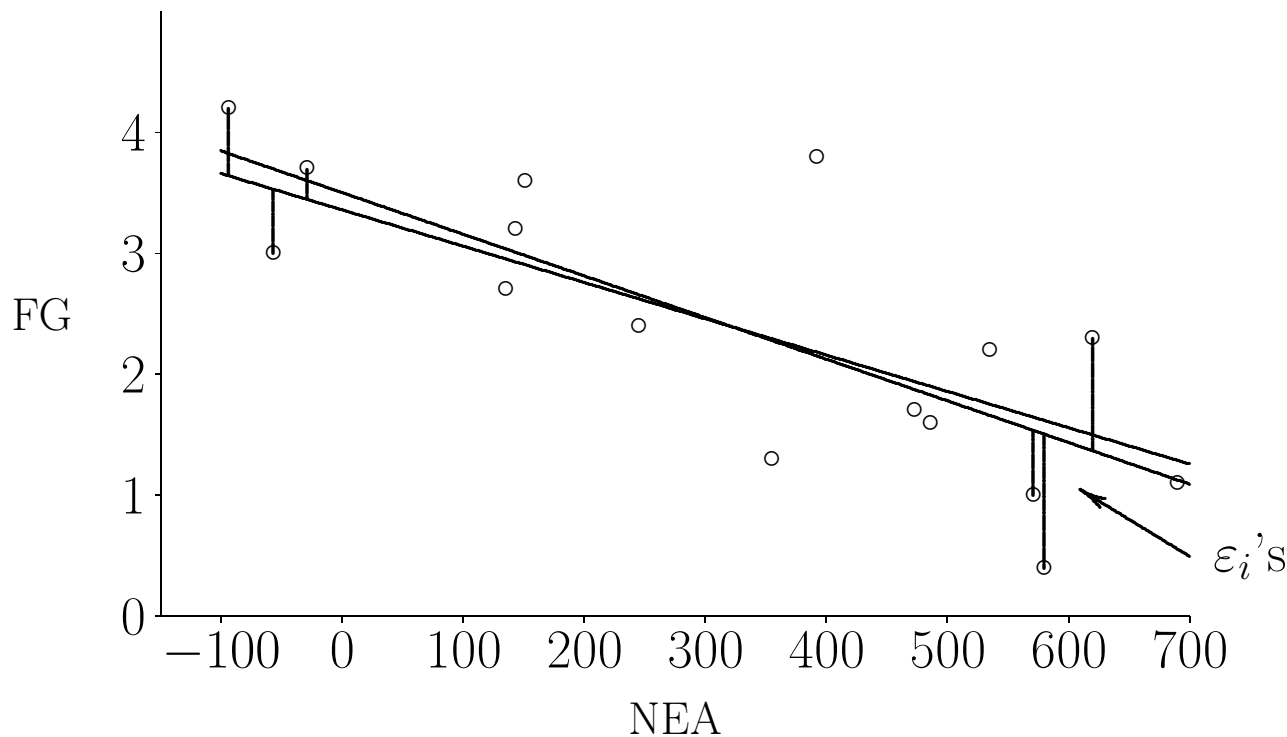
Statistical model:

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (= \alpha_0 + \beta_1(x_i - \bar{x}) + \varepsilon_i),$$

where the (*vertical*) errors $\varepsilon_1, \dots, \varepsilon_{16}$ are i.i.d. and $\sim N(0, \sigma)$,

- parameters: β_1, β_0 (or α_0) and σ ,
- x 's considered fixed — thus no capitals,
- assumptions:
 - * the linear relation: $EY_i = \beta_0 + \beta_1 x_i$,
 - * normal distribution of errors ε_i ,
 - * same stand.dev. of all observations (homogeneity),
 - * independence of errors (and of observations).

LEAST SQUARES ESTIMATION



How to choose the regression line (estimate β_0, β_1)?

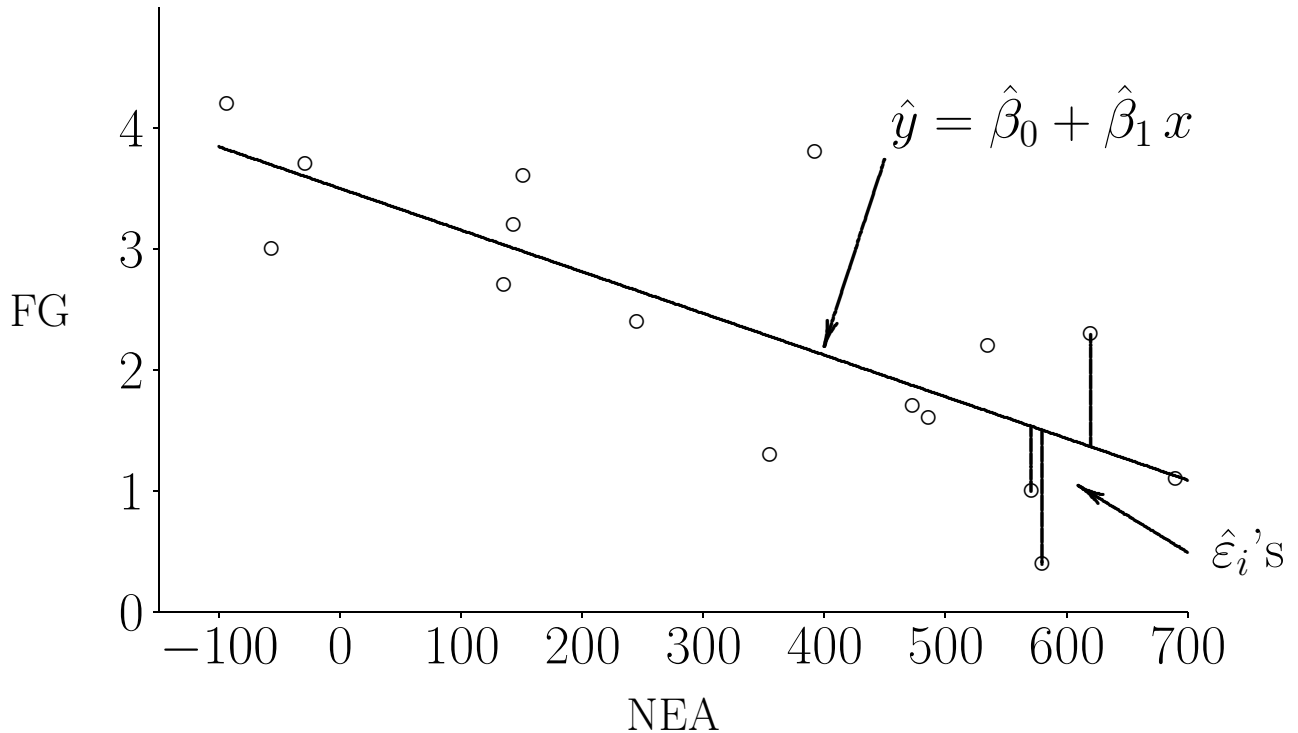
Idea: “best” line minimizes the sum of squared errors

$$\sum_i \varepsilon_i^2 = \sum_i (Y_i - \beta_0 - \beta_1 x_i)^2.$$

Motivations:

- intuitive (minimizes squared vertical deviations),
- easy to calculate (solutions have closed formulae),
- resulting estimates have good theoretical properties (unbiased and optimal for present model).

PARAMETER ESTIMATES

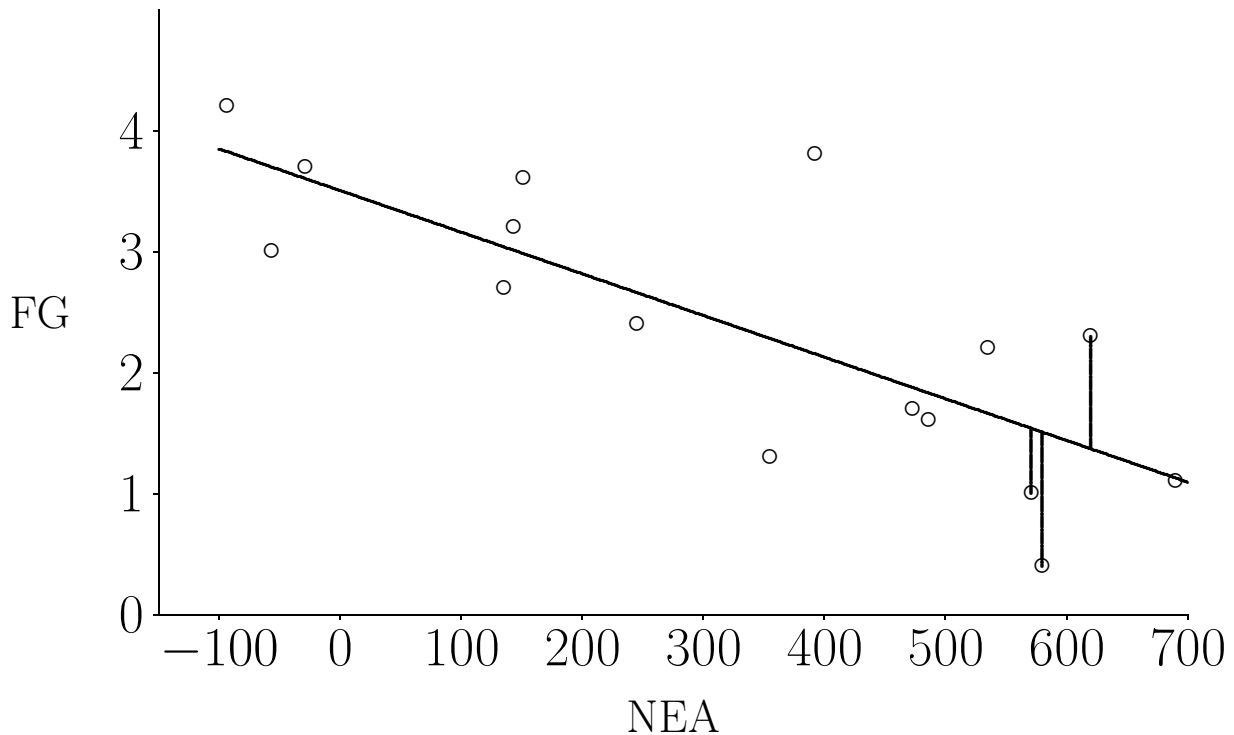


Parameter estimates:

- slope: $\hat{\beta}_1 = \frac{\sum_i (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} = r s_y / s_x$, ($r = \text{correlation}$),
- $\hat{\alpha}_0 = \bar{Y}$ (\Rightarrow estimated line passes through (\bar{x}, \bar{Y})),
- intercept: $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$,
- estimated line: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$,
- residual: $\hat{\epsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$, (“observed – predicted”)
- $\hat{\sigma}^2 = s^2 = \frac{1}{n-2} \sum_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \frac{1}{n-2} \sum_i \hat{\epsilon}_i^2$.

Minitab/Stata/R give the estimates and associated standard errors for mean parameters.

TESTS AND CONFIDENCE INTERVALS

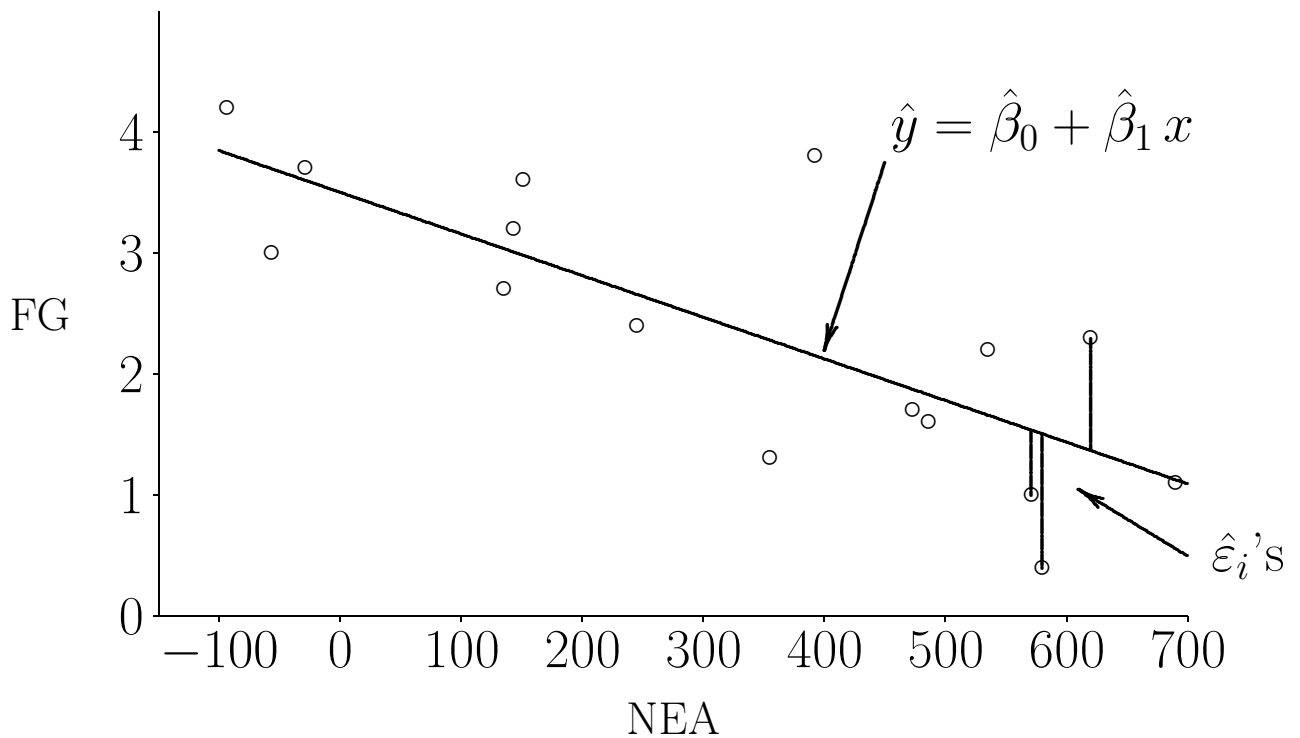


Statistical hypotheses about the parameters of the regression line are tested the “usual way”, using estimates and their standard errors:

- degrees of freedom for s^2 : $DFE = n - 2$,
- example: test of slope equal to known value:
 - * $H_0: \beta_1 = b$, (b known, fixed value),
 - * $H_a: \beta_1 \neq b$ (two-sided alternative),
 - * test: $t = (\hat{\beta}_1 - b) / SE(\hat{\beta}_1) \sim t(DFE)$ -dist. under H_0 ,
 most common example is $b = 0$ (horizontal line \sim no linear relation between x and y).

Confidence intervals: also the “usual way”, for example,
 95% CI: estimate $\pm t^* \cdot SE(\text{estimate})$, $t^* = t_{.975}(DFE)$.

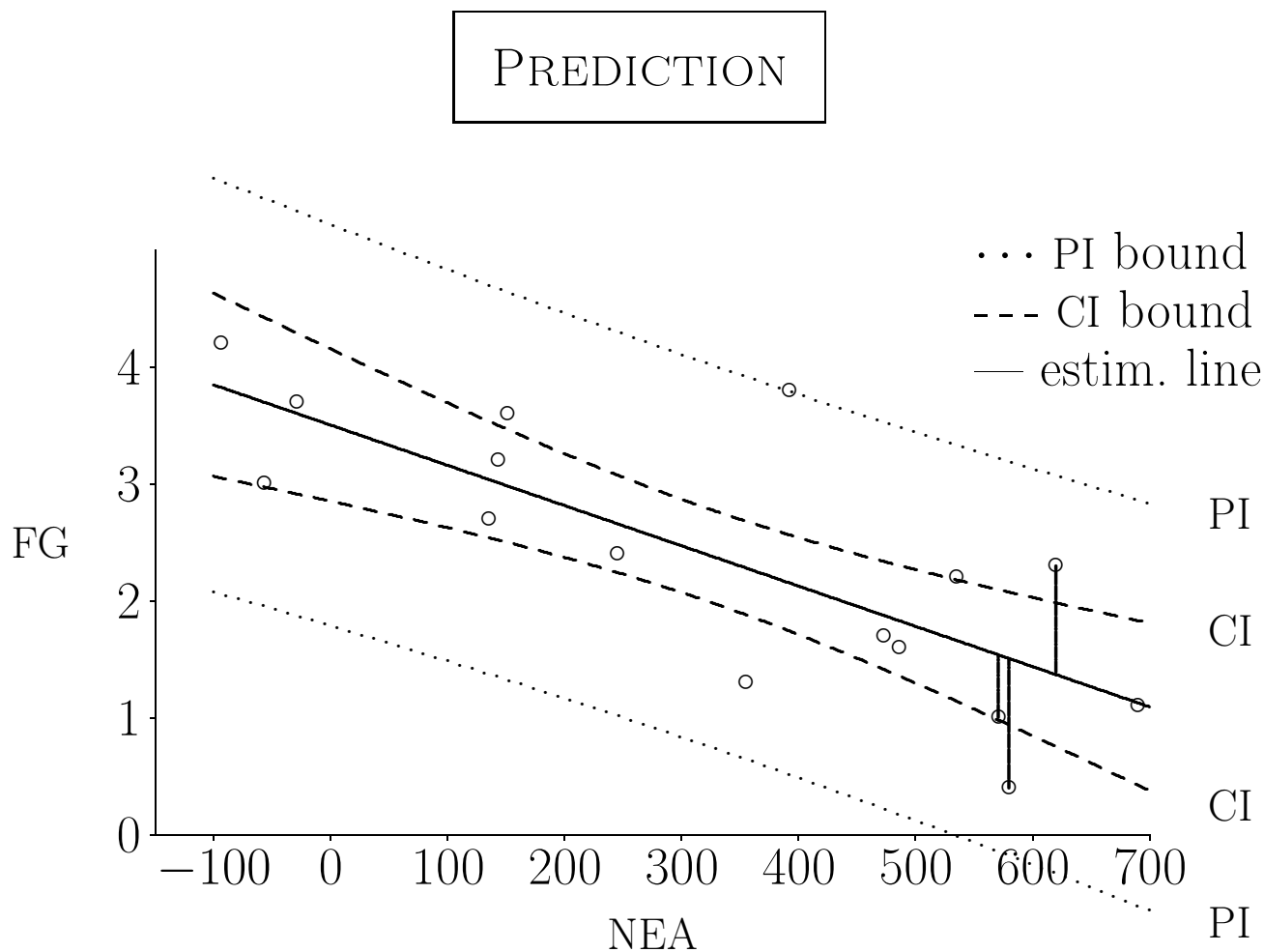
ANOVA TABLE FOR LINEAR REGRESSION



Hypothesis $H_0: \beta_1 = 0$ can also be tested by ANOVA:

Source of variation	Degrees of freedom	Sum of squares	Mean square	F
Regression Model	DFM = 1	SSM	MSM = SSM/DFM	MSM/MSE
Error	DFE = $n - 2$	SSE = $\sum_i \hat{\epsilon}_i^2$	MSE = SSE/DFE	
Total	DFT = $n - 1$	SST		

- estimated error variance = $s^2 = \text{MSE}$, as usual,
- F -test equivalent to t -test, because $F = t^2$, (same P)
- ANOVA table not really needed for simple linear regression (but for models with more x -variables).

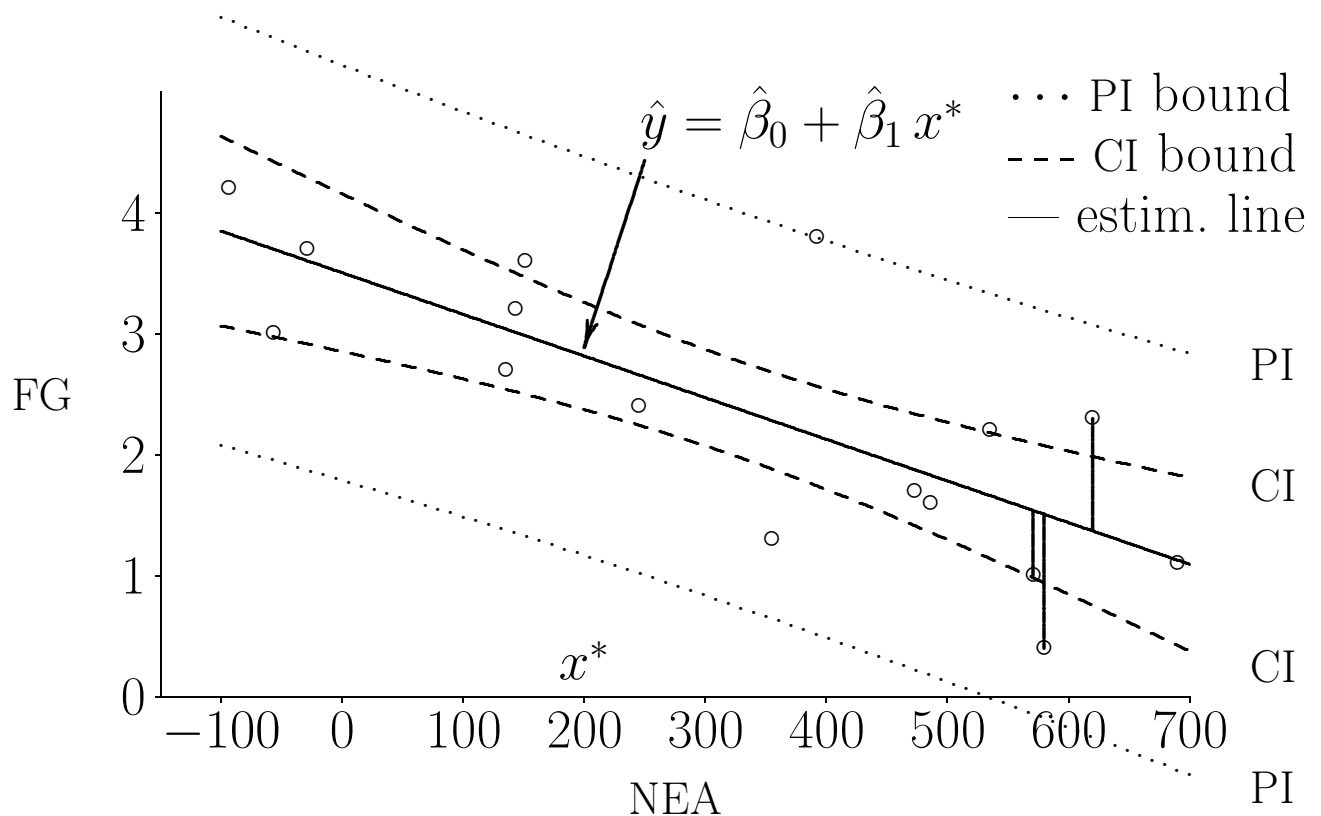


Prediction / Estimation:⁷

- 2 situations / purposes:
 - (CI) estimation of the regression line for given x , and CI to indicate the precision of the estimation,
 - (PI) prediction of a new observation for given x , and prediction interval (PI) to indicate *both* precision of the line (mean) and the dispersion around it,
- same estimated/predicted value: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$,
- CI for line more narrow than PI for new observation.

⁷ Stata terminology: prediction \sim estimation, forecasting \sim prediction.

STANDARD ERRORS IN LINEAR REGRESSION



Formulae for standard errors⁸ in linear regression:

$$\begin{aligned} \text{slope : } SE(\hat{\beta}_1) &= s / \sqrt{\sum_i (x_i - \bar{x})^2}, \\ \text{intercept : } SE(\hat{\beta}_0) &= s \sqrt{1/n + \bar{x}^2 / \sum_i (x_i - \bar{x})^2}, \\ \text{CI}(x^*) : SE(\hat{\mu}) &= s \sqrt{1/n + (x^* - \bar{x})^2 / \sum_i (x_i - \bar{x})^2}, \\ \text{PI}(x^*) : SE(\hat{y}) &= s \sqrt{1 + 1/n + (x^* - \bar{x})^2 / \sum_i (x_i - \bar{x})^2}. \end{aligned}$$

Notes:

- we usually don't compute by hand (calculator)!
- the squared variation of x 's determines accuracy,
- added "1" in prediction formula compared to estimation.

⁸ Strictly speaking, the error involved in a prediction is not a standard error but a prediction error; however, we stick here to the textbook notation/terminology.

HOW TO REPORT STATISTICS IN SCIENTIFIC PAPERS?

Reporting guidelines exist, both generally for statistics (Land & Altman, 2013) and for many specific study types (www.equator-network.org), and scientific journals increasingly require compliance with these.

Sample guidelines for medical journals (from Bailar & Mostellar, 1988⁹):

1. Describe statistical methods with enough detail to enable a knowledgeable reader with access to the original data to verify the reported results.
2. When possible, quantify findings and present them with appropriate indicators of measurement error or uncertainty (such as confidence intervals).
3. Avoid sole reliance on statistical hypothesis testing, such as the use of P values, which fails to convey important quantitative information.
5. Give details about randomization.
8. Give numbers of observations (and give the experimental unit).
9. Report losses to observation (such as dropouts from a clinical trial).
10. References for study design and statistical methods should be to standard works (with pages stated) when possible rather than to papers where designs or methods were originally reported.
11. Specify any general-use computer programs used.
12. Put general descriptions of statistical methods in the Methods section. When data are summarized in the Results section, specify the statistical methods used to analyze them.
13. Restrict tables and figures to those needed to explain the argument of the paper and to assess its support. Use graphs as an alternative to tables with many entries; do not duplicate data in graphs and tables.
14. Avoid non-technical uses of technical terms in statistics, such as "random" (which implies a randomizing device), "normal", "significant", "correlation", and "sample".
15. Define statistical terms, abbreviations, and most symbols.

⁹ See link at the VHM 801 homepage, people.ypei.ca/hstryhn/vhm801, for elaboration.

REVIEW: ANOTHER LOOK AT t AND t^*

From a (hypothetical) confused individual:

what are all those different t 's: t^* , $t(\text{df})$, $t_{1-\alpha}(\text{df})$, t_{obs} , t , ...

For confidence intervals, e.g. with *confidence level* 95% ($1 - \alpha$) and *error level* 5% (α), we need in our formula (1-sample),

$$\mu : \bar{X} \pm t^* s / \sqrt{n}, \quad t^* = t_{1-\alpha/2}(\text{df}),$$

a number, $t_{1-\alpha/2}(\text{df})$, from a t distribution with df degrees of freedom:

- it determines the middle 95% of the t distribution, between 2.5% and 97.5% ($\alpha/2$ and $1 - \alpha/2$),
- it is the 97.5% percentile in the t distribution,
- it can be found in statistical table (under confidence level 95% or tail area probability 2.5%),
- Minitab: enter the value 0.975 (or 0.025) and use **Inverse Cumulative Probability**,
- Stata: `invttail(df,0.025)`; R: `qt(0.025,df)`.

For a test of $H_0: \mu = \mu_0$ (where μ_0 is a known value), we compute the observed value of our t -statistic from the formula (1-sample),

$$t_{\text{obs}} = (\bar{X} - \mu_0) / (s / \sqrt{n}),$$

and the P -value is calculated from $P(t(\text{df}) > |t_{\text{obs}}|)$, where $t(\text{df})$ (or just t) refers to the t distribution with df degrees of freedom,

- it is a tail area probability, and can be evaluated as $<$ or $>$ specific values using table, based on table values below or above $|t_{\text{obs}}|$,
- Minitab: enter the value $|t_{\text{obs}}|$, use **Cumulative Probability**, and subtract the result from 1; or directly using the value $-|t_{\text{obs}}|$,
- Stata: `ttail(df,|t_obs|)`; R: `pt(-|t_obs|,df)`