Lecture 12 (Ch.3)

Last time we did Ordinary Least Squeres (OLS) regression, also known as curve/line fitting. It's a dangerously intuitive method. But at The 390 level, you are expected to know about a level of detail That is devilishly complicated. I.e. The devil is in The details. Look 1

In regression, we assume that data (xi, yi) follow a model: Model: 7 Yi= d+ B Xi + Ei evvor/residual. obs. y atx; y of line y(x) = x + Bx at x = x; To find The best "line (ie. DLS x, B), we minimized SSE: $SSE = \underbrace{\underbrace{\underbrace{\underbrace{\underbrace{x}}}_{i}}_{i} \in \underbrace{\underbrace{\underbrace{x}}_{i}}_{i} = \underbrace{\underbrace{\underbrace{\underbrace{x}}_{i}}_{i} \left[\underbrace{\underbrace{\underbrace{x}}_{i}}_{i} - \underbrace{\underbrace{a+\beta}_{i}}_{i}\right]^{2}$ Also, $\hat{\beta} = \frac{S_{xy}}{S_{xx}}$ where $S_{xx} = \sum_{i} (x_i - \overline{x})^2$ S_{xx} $S_{xy} = \sum_{i} (x_i - \overline{x}) (y_i - \overline{y})$ Compare [] When a, B are obtained from regression. Then, given x, we can predict y from The equ. of The "best" (ie.ons) fit: y (e.g. ICP, ie, Hard) $\frac{1}{\gamma(x) = x + \beta x}$ The "prediction", or "predicted value" or "fitted value", ... × × Blood Volocity (ie. easy)

Examples of BAD vegression



For new hatchlings and their mothers, noisier environments were associated with lower baseline levels of corticosterone, a hormone involved in stress responses. This negative correlation may seem counterintuitive, but a lower baseline has also been observed in humans with PTSD. Chronic stress can cause a sustained increase or decrease.

Most Likely to Succeed

In an ideal world, educators would know which students will best respond to which curricula. Researchers at Stanford University set out to see whether brain scans could help achieve this dream. They took 24 third graders and put them in a magnetic resonance imaging machine before they went through an eight-week one-on-one math tutoring program. Students whose scans showed a greater volume of tissue (gray matter) in the right hippocampus had higher performance gains from the tutoring than those with lesser volume in this brain area, which plays a critical role in forming new memories.



Scientific Anev. April 2018 Authors clain that leads of some hormone in the brain (Y-oxis) fall as the distance between bird nests and sources of sound increase. Ignore the regression line; would you agree with the conclusion, based on the scatter plot? What if we remove the 2003 circled birds?

Scientific Amer. March 2018

Authors claim The bigger a brain (x-axis) The more The gain from tutoring (y-axis). Ignore The regression line; would you agree with The conclusion, based on The scatter plot? What if we remove The lor 2 circled patients.

Both of These studies are highly sensitive to sampling. And even if The vegression lines are truly sloped, given The huge sprend in The scatterplots, do we care That The slopes are non-zero?



belo



Regression is a very powerful (sharp) tool. But, if you are not Caveful, you can do scribus damage !

I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

Complete change of perspective ! So far, we have viewed vegression as "curve-fitting." Pata: (xi, Yi) Model: $Y_i = \alpha + \beta x_i + \epsilon_i$ OLS estimates of $\alpha_i \beta$. Minimize: $SSE = \Xi \epsilon_i^2 \implies \hat{\alpha}, \hat{\beta} \implies predict: \hat{\gamma}(x) = \hat{\alpha} + \hat{\beta} x$ But regression can also be viewed as a method for error reduction. -> This other way of viewing regression will look very different, and it's also less intuitive/geometrical. But it is the more useful and generalizable view. In this view, the focus is on the variance of the y's, and as a result, the method is called The Analysis of Variance (ANOVA). It leads to quantities called R^2 and s_e which together quantify how good is the regression model. Let me motivate it: $S_{\gamma} = \sqrt{\frac{1}{n-1}} \frac{z}{i} \left(\gamma_i - \overline{\gamma}\right)^2 = 10 \text{ Cm}$ → Suppose we want to know a table's length. → We measure it, repeated by: Yi Say LOI times → make The histogram = 5yy = 10 cm Important (see next pag They page) So, if we use y as our estimate of The True table length, Then sy (or sign) gives us a measure of our uncertainty. And so, we can report Table Length: 150±10 cm (Y±sy) -> Now, suppose you are unhappy with the large sy. \sim -> You may wonder, could some of that variability be due to something else that is varying everytime you -> make a measurement of y- Eg. tengerature? - call This X. -> If so, then by measuring y and x, we may be able to reduce the ± of our report, by specifying y at a given x. Here is The math 2

The Analysis of Variance (ANOVA) Decomposition: Table length due to the (linear) How much of The variation in y is e temperature. relationship between y and r? Variance of $\gamma = \frac{1}{n-1} \frac{s_1}{\frac{s_1}{s_1}} \left(\frac{\gamma_1 - \gamma_1}{\gamma_1} \right)^2 \leftarrow Variance$ _ La's call This "Variation" See $l + \hat{l}_i - \hat{l}_i$, $\hat{l}_i = \hat{a} + \hat{\beta} \times i$ Error to $= \underbrace{S_{YY}}_{i} = \underbrace{\xi_{i}}_{i} \left(\underbrace{Y_{i}}_{i} - \underbrace{Y}_{i} \right)^{2} = \underbrace{\xi_{i}}_{i} \left(\underbrace{\widehat{Y}_{i}}_{i} - \underbrace{Y}_{i} \right)^{2} + \underbrace{\xi_{i}}_{i} \left(\underbrace{Y_{i}}_{i} - \widehat{Y}_{i} \right)^{2}$ Sseeplained SSunceplained Variation in y due to Variation in y (or explained by) 7. unexplained by? SStotal total variation unexplained by x in y. = SSeeplained + SSE a = something t something less Than 10⁴. SST es. 104 So, it we use if as our estimate of True length, Then The uncertaity is reduced from Syy (or SST) to SSE. > Sexplained is generally reported as a percentage of SST: R²= <u>SSeepl</u> (X100) = [percentage of The variability in Y temperature <u>SST</u> (X100) = That is due to (or can be explained by) R. (Bad model/fit) O < R² < 1 (Good model/fit) => Ssunerp, is reported as an "average", s² = SSE meaning of se = "typical error" $= \lim_{n \to 2} \frac{1}{i} \left(\frac{y_i - \hat{y}_i}{1 - \hat{y}_i} \right)^2$ funny Avg. ⇒ Altogether, for y (Table length): compare with Report y(x) ± Se error reduction, Sy=⊥ & (y, -7) ~ The ANOVA decomposition (above) assures Se < Sy Jargon: R² = Coeff. of determination Se = sample std. der. of errors

Last example); $\lim_{x \to \infty} (\gamma_{n} \times) \implies \widehat{\gamma}(x) = -755 + 13.28 \times$ Height (x) weight (y) $(\gamma - \dot{\gamma})$ Ŷ 200 - 1.5 200 **...** L 201.5 72 70 180 174.9 5.(65 120 (08.5 11.5 118 68 148.3 - 30.3 120 190 70 174.91 15.(70 72 68 7=161.6, Sy=39.5 regression, if I wanted to predict The next person's weight, I would say : Y ± Sy = 161.6 ± 39.5 (*) (see bottom) > with vegression: $SST = \sum_{i} (y_i - \overline{y})^2 = \cdots = 6251 \ll Varia bility in y is reduced$ from to $SSE = \underbrace{\sum_{i}}_{i} (\gamma_{i} - \widehat{\gamma}_{i})^{2} = (-1.5)^{2} + (5.1)^{2} + (11.5)^{2} + (-30.3)^{2} + (15.1)^{2} = 1307$ $R^{2} = \frac{SSexpl}{SST} = \frac{SST-SSE}{SST} = \frac{6251.2 - 1307}{6251.2} = 0.79.$ Meaning: 79% of The variability (or variation) in Y (weight, ...) is due to (can be explained by) X (height, ...). Cor Table length or --- or temperature, or . or temperature, or --- . The other piece of the decomposition: Se = 1307 = 21 pounds Meaning: The typical ervor (ie. deviation of y values from the fit) is about 21 pounds. Report: weight (or Table length): y ± 21 & with R=0.79 ICP -755+13,3(x) theight or FV *) Note That The initial uncertainty (39.5) has been reduced (21) by doing requession.

Picture for the ANOVA decomposition: depr. ervors = 4;-4;-~ SSE = Smaller Than 104 [SS unexp] $\sim S_{\gamma\gamma} = 10^4$ [SSTOTA] 50, when there is a (linear) relationship between & 4 Y, Then some portion of the variation in y can be attributed to (or explained by) x. That portion is SSeep., and the (unexplained) vest is SS unexp = SSE. So the variability in y (Sy) is reduced to SSE. When There is no relationship between x and y, Then The fig looks like below. Note that This situation is equivalent to the situation where we have data only on y, and Not on x stall. In That case the best prediction for every case is J (see hu): eurous = -1; - 4; = 4:-7 SST 104 D SSE 104 150 -In This case There is no veduction K in SST at all, as expected.

(FYI) If There is no & data, Then The OLS prediction is just I. Je. Y. What are The R² and Se? Nox. Yi=Y defn. of Sy². $S_{e}^{2} = \frac{SSE}{n-2} = \frac{Z_{e}(Y_{e}^{-}, \widehat{Y}_{e})^{2}}{n-2} = \frac{Z_{e}(Y_{e}^{-}, \widehat{Y}_{e})^{2}}{n-2} = (\frac{n-1}{n-2}) S_{Y}^{2} \implies S_{e} \sim S_{Y}$ $R^{2} = I - \frac{SSE}{SST} = I - \frac{\sum_{i} (Y_{i} - \widehat{Y}_{i})^{2}}{\sum_{i} (Y_{i} - \widehat{Y}_{i})^{2}} = I - \frac{\sum_{i} (Y_{i} - \widehat{Y}_{i})^{2}}{\sum_{i} (Y_{i} - \widehat{Y}_{i})^{2}} = 0 \implies R^{2} = 0.$ oo So, if we use I as our prediction, Then R=O (Bad), and Se~Sy, ie. The typical error ~ typical der. in y, ie. nothing gained. Another situation when nothing is gained is if we make vandom predictions, e.g. $\hat{\gamma}_i = random.$ Suppose The mean and The var. of These random predictions are The same as Those of observations ie. $\hat{q}_i = \text{Vandom}$ with $\bar{q} = \bar{q}$, $S_{\hat{q}} = S_{\hat{q}}$. The picture is But now, something strange happens: Although one can use The formula for R^2 to \bar{q} \bar{q} . have the usual interpretation (ie. percentage of var. in y, explained by a) because $\hat{y}_i = vandom$ are not OLS predictions. So, we don't have The AMOVA decomposition at all. Same objection applies to se. Again, The AMOVA decomposition is covert only for OLS 9; $\widehat{\gamma} = random are not OLS predictions.$ But The blue one has lower se. This doesn't contradict ArrovA, because -> Y The red is not OLS. Inshort, both have equal precision, but blue is more accurate

In-let 12 -1 a) Show That The sample mean of The predictions, i.e. $\hat{q} = \pm \xi \hat{q}_{\hat{e}}$ is equal to The sample mean of The y observation, i.e. \bar{q} . Hint: use $\hat{\alpha} = \bar{q} - \hat{\beta} \bar{x}$. b) Then, use part of to show That The sample mean of ervors, $\hat{\epsilon} = \pm \hat{\xi}_{i}(\gamma_{i} - \hat{\gamma}_{i})$, is equal to zero. Note: $\hat{\epsilon}_{i} = \gamma_{i} - \hat{\gamma}_{i}$. (FYI: This is why So is defined as (S. (Vi-Yi)), o Thermise, it's zero!) c) A property of 0L5 is $\sum_{i} \hat{e}_{i} x_{i} = (\chi_{i} - \hat{\chi}_{i}) x_{i} = 0$. You don't need to prove it; it's essentially $\sum_{y_{i}} s_{i} s_{i} = 0$ written differently. But do use that result to show that The correl. coeff. between The errors (\hat{E}) and The predictions $(\hat{\gamma})$ is zero. Hint: The numerator of The r of 4 and v is $\sum(u_i - \overline{u})(v_i - \overline{v})$. Also use the results of parts a and b. (hw/et12-2) Here are 3 important facts about OLS (Do not derive them): En prediction egn. $\hat{\gamma}_i = \hat{\lambda} + \hat{\beta} \, \chi_i$ $\Sigma \hat{\epsilon}_{i} \equiv \Sigma (Y_{i} - \hat{Y}_{i}) = 0$ < This is \$ SSE, witten differently. EEx;= E (Y;-Y;)x;= O - This is & SSE, written differently. Now, consider The following identity (That we used in lect): $\sum_{i=1}^{n} |Y_{i} - \overline{Y}_{i}|^{2} = \sum_{i=1}^{n} [(Y_{i} - \overline{Y}_{i}) + (Y_{i} - \overline{Y}_{i})]^{2}$ $= \underbrace{=}_{i} (\widehat{Y}_{i} - \underbrace{-}_{i})^{2} + \underbrace{=}_{i} (\underbrace{Y_{i} - \widehat{Y}_{i}})^{2} + 2\underbrace{=}_{i} (\widehat{Y}_{i} - \underbrace{-}_{i}) (\underbrace{Y_{i} - \widehat{Y}_{i}})$ Use The above 3 facts to show that The last term (cross-term) is the get used to This phrase. hw lect12 3: For the data shown here: x= 45, 58, 71, 71, 85, 98, 108 y = 3.20, 3.40, 3.47, 3.55, 3.60, 3.70, 3.80a) Compute the eq. of the OLS fit. b) Compute the total variation, SST. c) Decompose SST into explained and unexplained. d) Compute R2 and interpret it (in English), e) Compute the std. dev of errors, s e, and interpret it (in English). All by hand. You may use R to compute sums, means, std. deviations, but not a function that does regression or analysis of variance.