multiple regression (3.5 & 11.4) Lecture 14 (Ch. 3- cnd) So far, J predictor Simple regression: Y= x+ Bx + E polynomial regression: y= d+ Bix+B2 x2++++++++++++++ Both of these models fall under linear regression, because they are linear in the parameters alpha, beta 1, beta 2, This linearity is desirable, because it leads to a system of linear eqns (linear in the parameters alpha, beta 1, beta 2, ...), and so, it guarantees a unique soln. It is not restrictive because the model can be as nonlinear (in x) as you would like it to be. Even the next topic ("multiple regression") still falls under linear regression, and for the same reason. multiple predictors x, xz, ... Now, multiple linear regression (Several new concepts). $\pm g \cdot \gamma = \alpha + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_3 (x_1)^2 + \beta_4 (x_2)^3 + \beta_5 x_1 x_2 + \cdots$ E-g. "Vesponse" 2nd Variable/predictor, not z'd case. "Interaction y= Age at death, x_= income, x_2 = health Specific y = ICP, $x_i = blood$ flow, $x_2 = blood$ pressure. = vegression $y = \Delta Q$ (heat) $x_i = m(mass)$ $x_2 = \Delta T$ (temper.) $\Delta Q = C m \Delta T$ interaction T Pattern Recognition models are generally multiple regression Geometry: Instead of a line, we have a hyper-surface E.g. Y= d+ B, X, + B2X2 $\boldsymbol{\chi}_{\mathsf{I}}$

How to estimate &, B1, B2, --- Bx? Same as before, ic. with OLS => a, B, B, , --- Bn tow to do ANOVA? Same as before, because The decomposition depends on y's, not x's. The only "new" Thing is The count of The x's: SST = SSeepl. + SSumerplained $= \frac{1}{i} (\hat{\gamma}_i - \tilde{\gamma}_i)^2 = SSE.$ $R^{2} = \frac{SSerpl}{SST} = 1 - \frac{SSE}{SST}$ $P = \sqrt{\frac{SSE}{N - (k+1)}} = df$ Keep in mind That everytime you add a term on The R.H.S. One says that SSE has (ej a new predictor, df= n-(k+1). prosf, later. a non-linear term, an interaction, or even a completely random variable) k = # of ß's. you increase The chances of overfitting, E.g. ic. R² will increase (never decrease), $y = \alpha + \beta_1 x + \beta_2 x^2$, k + l = 3 $P_{ad_j}^2 = 1 - \frac{SSE / [n - (k+1)]}{SST / (n-1)} = 1 - \frac{Se^2}{S^2}$ $\gamma = \alpha + \beta_1 \chi + \beta_2 \chi^4$, =3 Reall R² > 1 as model gets more $\gamma = \alpha + \beta_1 x_1 + \beta_2 x_2 , = 3$ complicated. Radi attempts to fix that problem, but only $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 , = 4$ partially, I.e. both Rt and What determines The complexity of Ray never decrease as the a model is k+1, Not The nonlinearity model gets more complex. of The terms.

Read

Overfitting in maltiple vervessi

We know that one can overfit data on x and y, if one uses a high-order polynomial in polynomial regression. Recall that the main reason overfitting can happen is because such a regression model will have a lot of parameters, which in turn allow the fit to be more curvy/nonlinear, thereby going through every case.

In multiple regression there is yet another way that overfitting can happen even without high-order (nonlinear or polynomial) terms in the model.

Consider 3 cases on y and x_1. A model like $y = alpha + beta \times 1$, (a line) cannot over fit that data.

But a model like $y = alpha + beta_1 x_1 + beta_2 x_2$ (a plain) overfits completely. To see why, visualize the geometry of the problem; the reason for the complete overfit is that now the 3 cases are in 3D (not 2D), and there is always a plain that goes through any 3 points exactly.

%1

Note that the additional predictor x_2 can even be completely unrelated to y; it can even be just random numbers! In other words, by arbitrarily making the space big (by adding another predictor, even a useless one), we opened up the possibility of overfitting. So one can overfit even without any non-linear (e.g. quadratic, cubic, ...) terms.

You way think this is happening only because we are dealing with 3 cases here. But even with more cases, one can still overfit by simply including more (even random) predictors in the model. In general, if there are significantly more parameters/predictors than cases, then overfitting may happen. Furthermore, this overfitting problem is not specific to regression; All models can overfit ifwhen they are too large. Al/Machine Learning students, watch out!

Interpretation of B's

It looks like doing multiple regression is easy. And it is: $y = \alpha + \beta_1 \chi_1 + \beta_2 \chi_2 + \beta_3 \chi_1^2 + \beta_4 \chi_1 \chi_2 + \cdots$ $lm(\gamma \sim \chi_1 + \chi_2 + I(\chi_1^2) + I(\chi_1;\chi_2) + \cdots)$ But recall that regression is good not only for making predictions, but also for understanding what's going on in the problem. And for that, we interpret the B's. E.g. if y= x+2.3 x. - 2.1 xz, Then we can say that on The avg. y decreases by 2.1 units when xz increases by 1 unit. But there are 2 situations in which The B's connot be interpreted at all; They are associated with 2 concepts That we cover in this lecture : Collinearity and Interaction. Here is a real example: Data on y= life span $\gamma = \chi + 2.3 \chi - 2.1 \chi_2$ $x_i = health$ X2 = wealth o does y change +2.3 units when x, increases by 1 unit? · does y change -2.1 units when rez increases by 1 unit? so something is wrong with our interpretation of B! What?

Interaction

La's deal with interaction, first, because it's easier. $\frac{Q}{2} \quad \text{what does it look like?} \\ \gamma = \chi_1 \chi_2 \quad \longrightarrow \quad$ like XOR in logic y 1 x2 lines 7 what are the consequences of an interaction term? Q depends on other predictor(s) The effect of one predictor on y, A: - $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 =$ For This reason, The A's are G The effect of Az on y depends on x, un inter pretable, jie. They don't depends on x Tell you how much y changes when x is changed by 1 unit. (Guidance) Q: How do we know when to include an interaction term? A: 1) Look at scatterplots for signs of a saddle surface. 2) Include it, and see if it significantly improves R². 3) Ultimately, if it improves <u>predictions</u>, then include it! But always be mindful of overfitting. 4) Look at residual plots for signs of a saddle [skip This quarter] 5) In Ch.II we will leave another way of deciding whether or not an interaction exists,

Collinearity

Let's return to The first (important) step: Look at data ! multiple predictors => matrix of scatterplots: 12 Υ<u>΄</u> A linear association between χ_{1} x, and xz is called "collinearity" health It's a "bad" desease (See below, for why.) One cure is to NST. Collinearity simply exclude one of The predictors γz x1, x2 -Wedt lifespan (Almost any Thing is ok here. Don't look for a linear veloction. In fact, if you get a high-corr. linear pattern in one of Them, Then, you probably don't need The other one at all! See A consequence of collinearity is That it renders The B's un-interpretable (as The arg. rate of change of y ...): Ordinarily, in Y= x+ Rix, + Rzxz for But if x, and xz B1 = aug. rate of change in y, for 1 unit are correlated, The change in R., IF X2 IS HELD CONSTANT. One cannot hold one of them fixed In fast, in the lifespan = x + 2.3 (health) - 2.1 (wealth) example, The coeffs 2.3 and -2,1 cannot be interpreted as rates of change because health & wealth are linearly associated, i.e. There is collinearity. health wealth. buidance;) when you suspect extreme collinearity, Then just drop one of the predictors from the model. In less server cases, ask a statistican! Or look-up "principal component regression."

(FYE) when I collinearity Geometrically, the reason why the B's become uncertain and uninterpretated is that we are then trying to fit a plane through a cigar-shaped cloud in 3D, as opposed to a planar cloud. Collinearity ×1 ۲% د That is ambiguous! There are lots of planes one can fit Through a cigar-shaped cloud in 3D. Of course, Those different Fits differ in Their 2, B. Bz. That's why They become meaningless. You can also see That The predictions, y, are affected by collinearity; however, note that The effect is mostly in their uncertainty. (More, in Ch. 11). = Another bad consequence of coll. is That it effectively reduces The amount of information in The data, which, in turn, leads to more uncertain estimates of The B's and predictions. We'll see That in CR.II. = Another bad consequence of coll is That it can lead to evertitting. This is because the various predictors come with params to be estimated from data, but The varions predictors are essentially carrying The same information, i.e. There is effectively more params. Than data, hence overfitting can happen.

(vesidual plot FYI only \ One last Thing in regression (until Ch. 11): The visual assessment of fit quality. biased predictions (easy-to-fix) diagond N YX VS. Y : Random, symmetric about 6000 diagonal BAD < need a guodratic $\mathbf{\hat{s}}$ under/over forecasting small/large y Residual Plot: (Mona Lisa) ased predictions errors (last column of "table" before) better Ϋ́;-Ϋ́; Vi Dr (i) Not 0 D heed a quadratic YECO winended Random, symmetric under/over forecasting smill/large y. about error = 0 line BAD Model/fit] GOOD Model/fit things to do: so Nothing more to do 1) Transform Inta, or 2) fit diff. model (eg. polynomial or interaction--)



0.0

0

vesidual plot for Y= d+B> 1.0 1.5 2 yhat1

residual plot for y= + fix+ fix2

- yhat2 y=depx clearly shows That y=depx is The residual not a good fit to the data; because it shows a nonlinear pattern, we try y= d+px+ px2 and Then we see That The vesidual plot looks good (ie. shows no pattern). This may seen Trivial in This case, because you can see That you need y= x+ Bx+ fx2 from The scattenplot of The data, to begin with 1 But even in cases where you cannot make a scattenplot of data (e.g. in miltigle vegression) you can still make a residual plot.

In closing Ch. 3: best prediction ± uncertainty Before regression: 51 After repression: Se SSE $\widehat{\gamma} = \widehat{\alpha} + \widehat{\beta}_1 x_1 + \widehat{\beta}_2 x_2 + \cdots + \widehat{\beta}_3 x_1^3 + \cdots + \widehat{\beta}_4 x_1 x_2 + \cdots$ n-(k+1)df of SSE multiple predictors interaction k= # of B's in nonlinear tims terms Thanks to ANOVA, Se SSy (ie. regression reduces uncertainty)

Warning: Regression is a powerful tool, and so one can do serious damage with it. At the least, you must be aware of overfitting, interaction, and collinearity. There is a complex interplay between them, even though they are completely different concepts.

For example, even though both interaction and collinearity make the beta's un-interpretable, they are very different concepts: the former refers to a term in the regression model, or equivalently, a saddle surface in the 3d scatterplot. Either way, it involves x1, x2, and y. By contrast, the latter refers to (or is defined as) a linear relationship between the predictors x1 and x2, nothing to do with y. And yet, they both lead to uninterpretable beta's.

One may think that these 2 issues are not important if/when you are not interested in interpreting the beta's (e.g. if/when you are interested in making predictions only). And there is some truth to that thinking. But, these issues can affect your predictions, too! For example, both of them can lead to overfitting. Interaction can lead to overfitting because it introduces one more parameter into the model, allowing the model to twist and turn its way through more points. Collinearity can lead to overfitting because each of the predictors comes with its own beta in the model eqn, and yet, each of the two predictors is essentially carrying the same information. So, don't ignore the three concepts of overfitting, interaction, and collinearity.

Finally, and again, recall that multiple regression (even with polynomial terms), is still an example of linear regression, because the model y = ... is linear in the parameters. As a result of that linearity, when we take derivatives of SSE, we get a set of linear equations (again, linear in the parameters) that can be solved exactly, giving a unique solution. This is the advantage of linear models. Meanwhile, the models can be as nonlinear (in x) as you choose to make them by adding higher powers, or products, of the predictors. So, multiple regression is a win-win!

hw-let14-1 Consider fitting The model $Y_i = \alpha + \beta x_i; x_{2i} + \epsilon_i$, i = 1 - nto data on x_i, x_2, Y . Find The 2 normal equations of regression that one gets by differentiating $SSE = \sum_{i}^{\infty} (Y_i - \alpha - \beta x_i; x_{2i})^2$ with respect to α and β . Write The 2 equs in "bar" notation; but don't solve them for à, p.

hw lect14 2 (By R)

The article "The Undrained Strength of Some Thawed Permafrost Soils" (Canadian Geotech. J., 1979: 420-427) contained the accompanying data on y shear strength of sandy soil (kPa), x1 depth (m), and x2 water content (%). Obs Depth Content Strength

| 1 | <u> </u> | 215 | 147 | - | | | | | |
|-----|----------|------|------|---|--|--|--|--|--|
| 1 | 0.9 | 51.5 | 14./ | | | | | | |
| 2 | 36.6 | 27.0 | 48.0 | | | | | | |
| 3 | 36.8 | 25.9 | 25.6 | | | | | | |
| 4 | 6.1 | 39.1 | 10.0 | | | | | | |
| 5 | 6.9 | 39.2 | 16.0 | | | | | | |
| 6 | 6.9 | 38.3 | 16.8 | | | | | | |
| 7 | 7.3 | 33.9 | 20.7 | | | | | | |
| - 8 | 8.4 | 33.8 | 38.8 | | | | | | |
| ő | 6.7 | 07.0 | 16.0 | | | | | | |
| 9 | 6.5 | 27.9 | 16.9 | | | | | | |
| 10 |) 8.0 | 33.1 | 27.0 | | | | | | |
| 11 | 4.5 | 26.3 | 16.0 | | | | | | |
| 12 | 9.9 | 37.0 | 24.9 | | | | | | |

2.9 34.6 7.3 13

14 2.0 36.4 12.8

a) Perform regression to predict y from x1, x2, $x3 = x1^2$, $x4 = x2^2$, and $x5 = x1^2$; and write down the coefficients of the various terms.

b) Can you interpret the regression coeficients? Explain.

c) Compute R² and explain what it says about goodness-of-fit ("in English").

d) Compute s e, and interpret ("in English").

e) Produce the residual plot (residuals vs. *predicted* y), and explain what it suggests, if any.

f) Now perform regression to predict y from x1 and x2 only.

g) Compute R² and explain what it says about goodness-of-fit.

h) Compare the above two R² values. Does the comparison suggest that at least one of the higher-order terms in the regression eqn provides useful information about strength?

Skip parte)

i) Compute s e for the model in part f, and compare it to that in part d. What do you conclude?

hw lect14 3 (By R)

For each of the data sets a) hw_3_dat1.txt and b) hw_3_dat2.txt, find the "best" (OLS) fit, and report R-squared and the standard deviation of the errors. Do not use some ad hoc criterion (like maximum R2) to determine what is the "best" model. Instead, use your knowledge of regression to find the best model, and explain in words why you think you have the best model. Specifically, make sure you address 1) collinearity, 2) interaction, and 3) nonlinearity.

hw_lect14_4

In a problem involving x1, x2, y, suppose we have selected x1 such that it has a very (very) strong linear association with y. For simplicity, suppose the linear pattern "goes through" the origin. Similarly for x2, ie. it's highly correlated with y. On the adjacent diagram, and in the following order, draw the

81

72

a) 2d scatterplot in the x1-y plain.

b) 2d scatterplot in the x2-y plain.

c) 3d scatterplot, i.e. cloud of data in 3d. It may be hard to draw, but do your best in showing perspective.d) 2d scatterplot in x1-x2 plain.

Hint: This requires only 3d visualization of the type we did in the lecture!