

Announcements

- HW 2 regrade
- HW 3
- EXAMS graded

85+ Awesome

(75-85) Solid

(65-75) Passing

(60-65) Push!

< 60 Reevaluate Study Methods

Modeling Problems in Regression

Sometimes we encounter issues which can harm the validity of our model fits or their inference.

Summary:

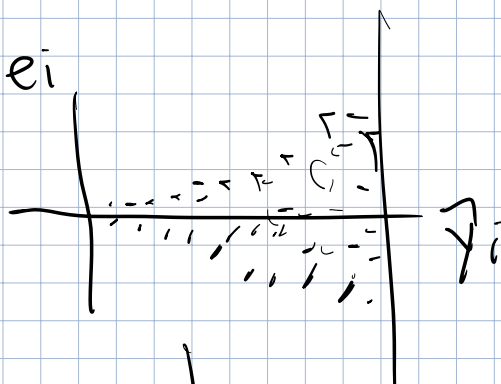
— Structural Problems:

- Multicollinearity
- Influential points

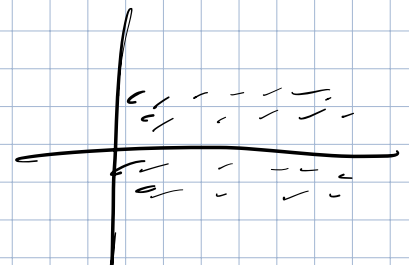
— Violations of Model Assumptions

- Heteroskedasticity
- Non-Normal residuals
- False assumption of linearity

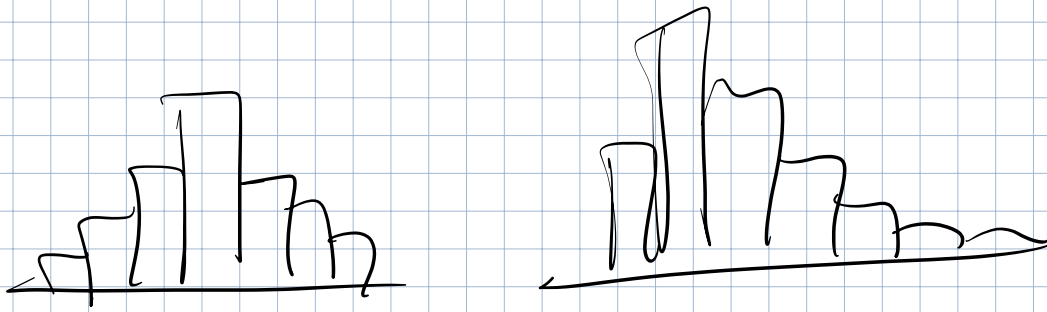
} next
wk



$$\sigma^2(x) = \sigma^2 x$$



↓ vs



For each problem, I'll try to give you:

- A. Damage caused
- B. Defection
- C. Solution

Let's start w/ Multicollinearity

Problem: Two or more of the predictors are highly correlated.

X ← two T of columns are (almost)
linearly dependent

\Rightarrow deficient rank $\text{rank}(X) < p$

$$\Rightarrow \hat{\beta} = \underbrace{(X^T X)^{-1}}_{\uparrow} X^T Y$$

may not exist, or is highly unstable

What do I mean by unstable?

$$X^T X = U \Lambda U^T$$

U = orthogonal matrices

Λ = diagonal matrix

$$(X^T X)^{-1} = (U \Lambda U^T)^{-1} \quad \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & 0 \end{pmatrix}$$

$$= (U^T)^{-1} \Lambda^{-1} U^{-1}$$

$$= U \Lambda^{-1} U^T$$

$$\uparrow = \begin{pmatrix} 1/\lambda_1 & & & \\ & 1/\lambda_2 & & \\ & & \ddots & \\ & & & 1/\lambda_n \end{pmatrix}$$

$$\text{var}(\hat{\beta}) = \sigma^2 \underline{\underline{(X^T X)^{-1}}}$$

When multicollinearity is present:

- the coefficients can swing wildly based on which other predictors are in the model
- the coefficient is highly sensitive to small changes in the model (in the X matrix)

Symptoms

Adding or dropping predictors causes massive change in coeffs & maybe even sign flips.

EX:

$\text{corr}(X_1, X_2) = 0$
A dataset w/ no multicollinearity

Model	$\hat{\beta}_1$	$\hat{\beta}_2$
$Y \sim X_1$	-1	NA
$Y \sim X_1 + X_2$	-1	-2

↳ no major change

$\text{corr}(Z_1, Z_2) = -9$
A dataset w/ severe multicollinearity

Model	$\hat{\beta}_1$	$\hat{\beta}_2$
$Y \sim Z_1$	-1	NA
$Y \sim Z_1 + Z_2$	5	-3

↳ major change

1
 $\Rightarrow se(\hat{\beta}_j)$ are getting inflated

Effect of this:

① when we do a t-test:

$$t = \frac{\hat{\beta}_j - 0}{\hat{se}(\hat{\beta}_j)} \uparrow \quad \downarrow$$

the inflation in $se(\hat{\beta}_j)$ makes it harder to reject $H_0 \Rightarrow$ to detect a signal when you have many related predictors in the model.

② If you look @ `anova(typ=1)` & change the order, the significance / t-stat / p-val may change dramatically

Note:

— In reality, some degree of multicollinearity always exists. Our job is to judge how much we are willing to live with.

- It becomes hard to understand the effects of individual predictors if they're all highly correlated.

Detection Methods

- ① Naive \rightarrow look @ correlation matrix of predictor

Watch out for very high r values in magnitude b/w predictors

- ② VIF = variance inflation factor

VIF measures how much the variance of the β_j s is inflated by adding that specific pred to a pre-existing model:

for the j^{th} predictor:

$$VIF_j = \frac{1}{1 - R_j^2} \quad \text{where}$$

R_j^2 is the coef. of determination when

you regress $X_j \sim X_1 + X_2 + \dots + X_{j-1} + X_{j+1} + \dots + X_{p-1}$

- If X_j is uncorrelated w/ all other preds. -

$$VIF = 1$$

- If $1 \leq VIF \leq 4 \Rightarrow$ "light multicollinearity"

$4 \leq VIF \leq 10 \Rightarrow$ "moderate"

$VIF > 10 \Rightarrow$ "severe"