# Modeling Problems

## Structural Problems
- Multicollinearity
- Influential Pts

## Violations of Model Assumptions
- Heteroskedasticity
- Non-Normal residuals
- False assumption of linearity

# Multicollinearity

$\boxed{\text{Problem}}$ two or more predictors are highly correlated

$$\Downarrow$$

could cause numerical difficulties when calculating

Design Matrix

$$X = \begin{pmatrix} \mathbb{1}_n & X_1 & X_2 & \cdots & X_{p-1} \end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad \longleftarrow \quad \text{LS Estimates}$$

$$\text{rank}(X) = p$$

For ex:

$$X_1 \qquad X_2 = \underline{4X_1}$$

$$y_i = 1 + \underline{2}X_{1i} + \underline{4}X_{2i} + \varepsilon_i \qquad \vec{\beta}_2 \begin{pmatrix} 1 \\ 2 \\ 4 \end{pmatrix}$$

$$\Downarrow$$

$$y_i = 1 + 2X_{1i} + 4(4X_{1i}) + \varepsilon_i \qquad \vec{\beta}^* = \begin{pmatrix} 1 \\ 18 \\ 0 \end{pmatrix}$$

$$= 1 + \underline{18}X_{1i} + \varepsilon_i + \underline{0}X_{2i} +$$

$$\Downarrow$$

$$y_i = 1 + 22X_{1i} + (-1)X_{2i} + \varepsilon_i$$

non-identifiability

# Damage

If I don't address this issue
here's some negative consequences:

① $\hat{\beta} = (X^TX)^{-1}X^Ty$

$$X^TX = U\Lambda U^T \qquad \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \approx 0 \end{pmatrix}$$

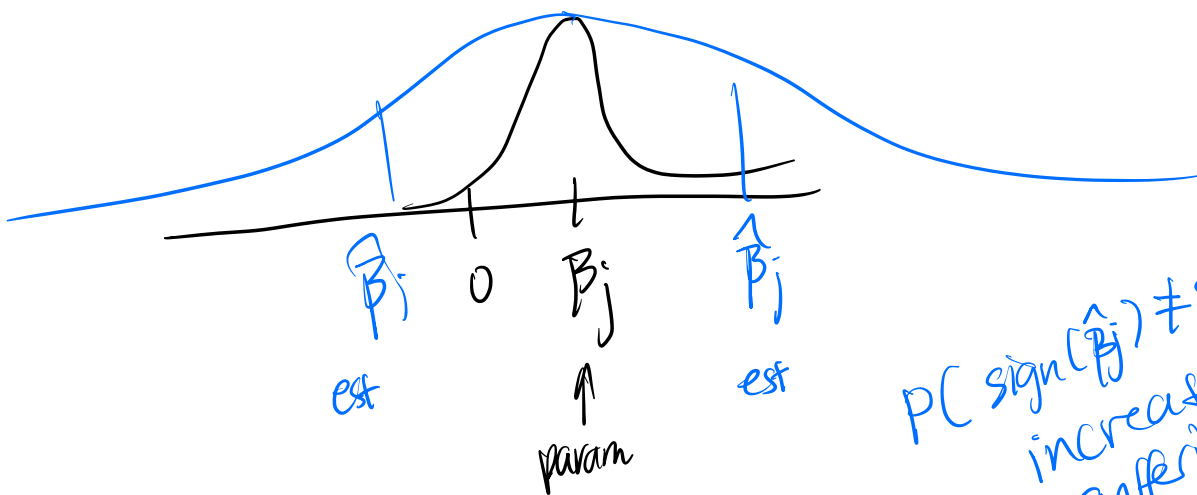$$(X^TX)^{-1} = U\Lambda^{-1}U^T \qquad \Lambda = \begin{pmatrix} 1/\lambda_1 & & \\ & 1/\lambda_2 & \\ & & \ddots \\ & & & 1/\approx 0 \end{pmatrix}$$

② $Var(\hat{\beta}) = \sigma^2 (X^TX)^{-1}$

unaddressed multicol. leads to inflated
standard errors.

# Impact on Inference:

③ $$t = \frac{\hat{\beta}_j - 0}{\hat{SE}(\hat{\beta}_j) \uparrow}$$

$\downarrow$ losing statistical power



$\hat{\beta}_j$ est    0    $\beta_j$ param    $\hat{\beta}_j$ est

$P(\text{sign}(\hat{\beta}_j) \neq \text{sign}(\beta_j))$ increased when suffering from multicollinearity

## Symptoms to look for:

— When you add a predictor to the model, the coefficients swing wildly, sometimes even changing signs.

— this is b/c the $\hat{\beta}$ is super sensitive to small changes in $X$.
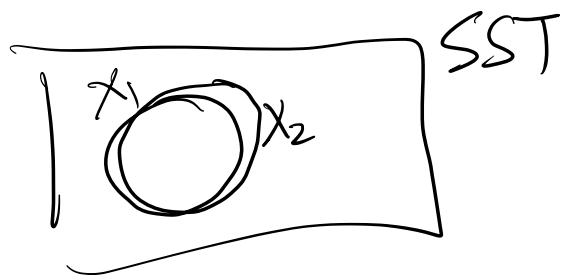
EX:

**NO multicol**
$Corr(X_1, X_2) = 0$

| Model | $\hat{\beta_1}$ | $\hat{\beta_2}$ |
|---|---|---|
| $Y \sim X_1$ | $-1$ | NA |
| $Y \sim X_1 + X_2$ | $-1$ | $-5$ |

**severe multicol**
$Corr(X_1, X_2) = .9$

| Model | $\hat{\beta_1}$ | $\hat{\beta_2}$ |
|---|---|---|
| $Y \sim X_1$ | $-1$ | NA |
| $Y \sim X_1 + X_2$ | $10$ | $-5$ |

## Symptom

ANOVA (typ=1) $\longrightarrow$ significance may change depending on order


SST

| | SS | F | P | |
|---|---|---|---|---|
| $X_1$ | ↑ | ↑ | ↓ | ✗ |
| $X_2$ | ↓ | ↓ | ↑ | NS |

Compare vs.

|       | SS | F | p |   |
|-------|----|----|----|----|
| $X_2$ | ↑ | ↑ | ↓ | * |
| $X_1$ | ↓ | ↓ | ↑ | NS |

## Notes:

- In reality multicollinearity is somewhat always present. Our job is mainly to decide how much we are ok with.

- Unchecked multicollinearity makes it very hard to understand the effect of each predictor on the response.

## Detection

① Correlation Matrix    (Naive)

② VIF — Variance Inflation Factor

VIF measures how much the variance of $\hat{\beta}$ are inflated by adding a specific predictor to the model

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

$\boxed{R_j^2}$ is the coef. of determination when we regress

$$X_j \sim X_1 + X_2 + \cdots + X_{j-1} + X_{j+1} + \cdots + X_{p-1}$$

If $\text{VIF} = 1 \iff$ no correlation b/w $X_j$ & other preds

$1 \leq \text{VIF} \leq 4 \iff$ "light" multicol.

$4 \leq \text{VIF} \leq 10 \iff$ "moderate"

$\text{VIF} \geq 10 \iff$ "severe"

Solutions

① Drop some of the suspicious looking predictors based on VIF

② Feature engineer the highly correlated variables into a single new predictor which

carries summarizes the info in the original preds.

More Complizated Approaches:

① Regularized Regression
$\left\{\begin{array}{l}\text{Ridge} \quad \nearrow \quad ||\beta||_2^2 \text{ penalty} \\ \\ \text{LASSO} \quad \nwarrow \quad ||\beta||_1 \text{ penalty}\end{array}\right.$

② Dimension Reduction on X
eg. PCA

③ Partial Least Squares