

Basics

$0 \leq \Pr(A) \leq 1$ sum of probs of all events in sample size is 1

$$\Pr(A) = 1 - \Pr(A^c)$$

Joint prob: $\Pr(A \cap B)$

$$\text{either } A \text{ or } B: \Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$$

Mutually exclusive: $\Pr(A \cap B) = 0$

Bayes

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

sensitivity

$$P(D=1 | R=1) = \frac{TP}{TP + FN}$$

maximize sensitivity
want to minimize false negs
ex: detecting cancer

specificity

$$P(D=0 | R=0) = \frac{TN}{TN + FP}$$

maximize specificity
want to minimize false positives
ex: costly tests

hw example

5322 male non smokers

7019 male smokers

↳ 16 non develop

↳ 71 smokers develop

		develop	don't develop	total
		smoker	non smoker	
smoker		71	6942	7019
non smoker		16	5306	5322
total		90	12248	12341

$$P(\text{non smoker develops}) = \frac{16}{5322}$$

$$P(\text{developing}) = \frac{90}{12341}$$

X and Y independent if $P(X|Y) = P(X)$ for all X, Y

or $P(Y|X) = P(Y)$

$$P(X \cap Y) = P(X) \cdot P(Y) \text{ for all } X, Y$$

ex: pop $\mu = 100 \quad \sigma = 20 \quad n = 64$

$$P(\bar{X} > 105)$$

$$P(\bar{X} > 105) = P\left(\frac{1}{n} \sum X_i > 105\right)$$

$$\frac{\bar{X} - 100}{\frac{20}{\sqrt{64}}} \sim N(0, 1)$$

$$P\left(\frac{\bar{X} - 100}{\frac{20}{\sqrt{64}}} > \frac{105 - 100}{\frac{20}{\sqrt{64}}}\right)$$

$$P\left(\bar{X} > \frac{5}{\frac{20}{\sqrt{64}}}\right)$$

computing variance

$$\text{var}(X) = E(X^2) - E(X)^2$$

$$E(X^2) = \sum x^2 P(x=x)$$

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$$

confidence intervals

wrong interpretations

- 95% prob that true param lies within the interval
why incorrect? the population parameter is a constant, it either falls in the interval or doesn't
 - 95% of all data in the population fall in the interval

as CI% increases, α decreases so $Z_{\alpha/2}$ gets larger

- 90% CI $\rightarrow Z \approx 1.645$
- 95% CI $\rightarrow Z \approx 1.96$ \Rightarrow so more confidence
- 99% CI $\rightarrow Z \approx 2.576$ = wider interval
= less precision

Correct

- if you repeat sample process many times , about 95% of the computed intervals will capture true pop. parameter
then also 5% of them will not capture pop. mean

MLE example ①

$$\text{Bernoulli} : P^{x_i} (1-P)^{1-x_i}$$

Given x_1, x_2, \dots, x_{100}

$$L(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$l(p) = \sum_{i=1}^n x_i \log(p) + (1-x_i) \log(1-p)$$

$$\frac{\partial}{\partial p} l = \sum_{i=1}^n \frac{x_i}{p} - \frac{1-x_i}{1-p}$$

$$S_{\text{et}} = 0 \Rightarrow \frac{1}{P} \sum_{i=1}^n x_i = \frac{1}{1-P} \sum_{i=1}^n (1-x_i)$$

$$(1-p) \sum x_i = p \sum (1-x_i)$$

$$(\sum x_i) - p \underbrace{\sum x_i}_{\text{in } S} = np - p \underbrace{\sum x_i}_{\text{in } S}$$

$$\sum x_i = np$$

13

11

$$\text{bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

SE : standard dev of estimator or sampling distribution

$$\text{Var}(\hat{\theta}) = E(\hat{\theta} - E(\hat{\theta}))^2$$

$$MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2 = SE^2 + Bias^2$$

consistency: estimate converges to true value as $n \rightarrow \infty$

$$\nearrow \text{SE} = \frac{\sigma}{\sqrt{n}}$$

C.I.T: x_1, x_2, \dots, x_n have pop mean μ , var σ^2

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

regardless of distribution of x_1, \dots, x_n

Hypothesis Testing

NUL (H₀)

ALTERNATIVE (H_A)

FRAMEWORK

- start with two hypotheses
- choose sample
- assess how likely sample data can be obs. given H₀ true
- if data are very unlikely then reject H₀
- otherwise fail to reject

never accept null *

P value

prob of obs. data or even more extreme values assuming H₀ true

$$H_0: \mu = 0; H_A: \mu > 0$$

$$P \text{ value} = P(X > X_{\text{obs}}) \text{ when } \mu = 0$$

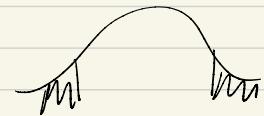
$$\text{if } H_0: \mu = 0, H_A: \mu \neq 0$$

$$P \text{ value} = 2 \times P(X > X_{\text{obs}}) \text{ two sided}$$



Significance level

predetermined threshold to see if result is stat. significant



if $P < \alpha = 0.05 \Rightarrow$ reject H₀

Approaches

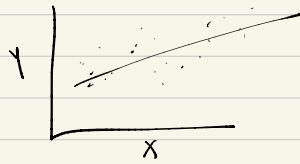
① simulation based \hookrightarrow simulate / resampling NUL

② parametric

\hookrightarrow assumptions of pop distib.

\hookrightarrow use CLT

Simple Linear Regression



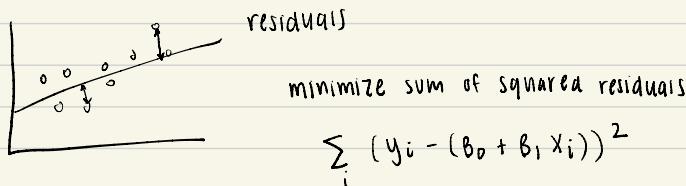
line of best fit between two continuous variables

$$y = \beta_0 + \beta_1 x + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$$y \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

↑ ↑ ↑ ↑
random variable parameter fixed variable parameter

use least squares estimation to find best β_0, β_1



$$\text{regression line } \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

residuals $r(\hat{e})$ as $r = y - \hat{y}$

Hypothesis Testing for coefficient estimates

$$H_0: \beta_1 = 0$$

$$\frac{\hat{\beta}_1}{SE(\hat{\beta}_1)} \sim t_{n-2}$$

$$y = 0.0498x + 14.473$$

MULTIPLE Linear Regression

why? most relationships cannot be fully explained by a single predictor and outcome

- many outcomes depend on multiple factors acting simultaneously

Confounding variables

confounder is

- ① correlated with predictor
and
- ② influences outcome independently

→ include confounders in regression model to help isolate true effect of each predictor

Directed Acyclic Graph

- show potential causal relationships

nodes = vars

arrows = causal effects

acyclic = no loops

help understand which vars to use as control & which cause bias

Multiple Linear Regression Model

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma^2)$$

y_i : outcome for obs i

assumed i.i.d

x_{ij} : val of predictor j for obs i

β_0 : intercept = expected value of y when all predictors are 0

β_j : expected change in y for a one unit change in x_j holding all other predictors constant

$$y_i \stackrel{\text{iid}}{\sim} N(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}, \sigma^2)$$

$$E[Y | X_1 = x_1, \dots, X_p = x_p] = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

interpretation different if predictors are categorical vars not continuous

Which variable is the strongest predictor?

use t-statistic $\Rightarrow t_j = \frac{\hat{\beta}_j}{\text{SE}(\hat{\beta}_j)}$

abs val.
larger $|t|$: stronger evidence that predictor is associated

= coeff. estimate over corresponding SE

coeffs are sensitive to scale of predictors but T is not
so magnitude of β alone is not reliable

Design Matrix

rows = # obs or sample size

cols = # parameters or predictors (include intercept *)

$$Y = X\beta + \varepsilon$$

$\uparrow \quad \uparrow \quad \uparrow$
 $n \times 1$ vector of random errors
 $n \times (p+1)$ matrix of responses
 $(p+1) \times 1$ vector of coefficients
 β

F test : evaluate whether **overall** model is significant

is there a relationship between the predictors and response?

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$$

None of the predictors explain variation in Y

H_A : at least one predictor is linearly associated with Y

if p value small (< 0.05) reject H_0

Compares variation explained by model & unexplained variation

if model explains more variation than random chance \rightarrow large F statistic

$$F \text{ statistic} = \frac{\text{mean square regression (MSR)}}{\text{mean square error (MSE)}}$$

Nested F Test (Type III test)

may want to assess association b/w categorical var & outcome

if there are dummy vars, test if a subset of coeffs = 0

- compares reduced model to full model

example

x_1 : age

x_2, x_3 : dummy vars for region

$$\text{full model: } Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$$H_0: \beta_2 = \beta_3 = 0$$

subset of region only coeffs

reduced model: says what if null is true $\Rightarrow Y = \beta_0 + \beta_1 x_1$

Nested F test asks: do additional predictors x_2 and x_3 significantly improve model's fit

t-test: test individual coefficients

for each coefficient β_j : test $H_0: \beta_j = 0$ vs $H_A: \beta_j \neq 0$

(given all the other vars in the model)

$$\text{test statistic } T = \frac{\hat{\beta}_j - 0}{\text{SE}(\hat{\beta}_j)} = \frac{\text{estimate} - NVII}{\text{SE}}$$

$$\text{confidence interval CI} = \hat{\beta}_j \pm \text{SE}(\hat{\beta}_j) C_\alpha$$

$$\curvearrowright C_\alpha = t_{\alpha/2, df}$$

Estimating Linear Regression Coeffs

Ordinary Least Squares

- OLS estimate obtained by minimizing residual sum of squares (RSS)

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - [B_0 + B_1 x_{i1} + \dots + B_p x_{ip}])^2$$

find OLS estimates by
minimizing by taking
partial derivs wrt each B_j

OLS finds vals of B_0, B_1, \dots, B_p to minimize this fn

Key facts

- no distributional assumptions about y or ϵ_i for estimation
- unbiased, min-variance estimates
- method of estimation, not inference
 - provides best fitting line through data by minimizing squared residuals
 - limited to linear regression

Multiple linear regression model

$$y_i = B_0 + B_1 x_{i1} + \dots + B_p x_{ip} + \epsilon_i \quad \begin{matrix} \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \\ i = 1, \dots, n \end{matrix}$$

$$\text{matrix representation } Y = X B + \epsilon \quad \epsilon \sim N(0, \sigma^2 I)$$

$$\text{OLS estimates } \hat{B} = (X^T X)^{-1} X^T Y$$

Properties of OLS estimator

$$\text{linearity } Y = X B + \epsilon$$

$$\text{Unbiasedness } E[\hat{B}] = B$$

$$\text{variance } \text{Var}(\hat{B}) = \sigma^2 (X^T X)^{-1}$$

* OLS gives the smallest variance among all linear unbiased estimators *

Categorical Variable Terms

levels = values of categorical var.

binary var = two levels

factor (specific to R) = categorical var that stores levels & labels

dummy var = numeric var that represents a categor. var (0/1)

reference / baseline level = value to compare other values of categorical var to
(important for interpreting coeffs)
for categorical dummies, reference = omitted category

If a categorical var has K levels, we need K-1 dummy vars

↳ including all K dummy vars causes perfect multicollinearity \Rightarrow sums to 1

interpreting categorical predictors : each coeff = mean diff between a group & reference category

↳ always interpret in relation to reference category

example

$$x_i = \begin{cases} 1 \\ 0 \end{cases}$$

example: avg diff in dependent var btwn specific category & reference category
ex/ pos coeff for female with male as reference, where predicting y=income
on avg, female predicted to have coeff more income than male
if there are multiple dummies for categorical, coeff is diff btwn the dummy & reference,
holding all other vars & dummies constant

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i \\ \beta_1 + \epsilon_i \end{cases}$$

Nested F Test / type III test

- want to assess association between categorical var & outcome
- since there are dummy variables test if subset of coeffs = 0

idea: test whether categorical var as a whole (with dummy vars) is signif. related to outcome

HOW?

- compare a reduced model (without categorical var) to full model (with included)

$$\text{Null: } H_0 = \beta_{j1} = \beta_{j2} = \dots = \beta_{jk} = 0$$

i.e. testing whether all dummy vars for a factor are zero

* USE ANOVA IN R?

interpretation: if p value small \rightarrow categorical var (as a group) significantly improves the model

Interaction Terms

how the relationship between a predictor and Y changes based on another (categorical) predictor

i.e. how x_1 and Y relationship changes depending on x_2

ex: how drug dosage & anxiety level relationship changes for those < 65 yrs v.s. ≥ 65 yrs

model representation

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 (X_1 * X_2) + \epsilon$$

$\nearrow \quad \nearrow \quad \underbrace{\quad \quad \quad}_{\begin{array}{l} \text{main effects} \\ (\text{each var w/o interaction}) \end{array}} \quad \underbrace{\quad \quad \quad}_{\begin{array}{l} \text{extra effect} \\ \text{of changing both} \\ \text{simultaneously} \end{array}}$

*but can be diff. to interpret main effects when interaction is signif.

β_3 = additional change in slope of X_1 for each unit increase in X_2

If interaction term (β_3) is statistically significant, the effect of one var depends on the value (or category) of the other

↳ can have higher order interactions ($X_1 X_2 X_3$) or continuous var interactions but hard to interpret

When to include interaction terms?

- theory/prior domain research suggests vars interact
 - EDA reveals diff slopes or group specific patterns
- ↳ don't add interactions blindly → increase model complexity & can reduce interpretability

Assessment

mean square error

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

R square proportion of variance

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

← residual sum of squares
← total sum of squares

Ratio of explained variation to total variation

R^2 is now well model explains variability in response var

$$R^2 \in [0, 1]$$

if $R^2 = 0 \Rightarrow$ model explains none of the variation

if $R^2 = 1 \Rightarrow$ model perfectly fits data

R^2 increases with additional predictors (even irrelevant ones)

Residual sum of squares

RSS = total squared diff btwn actual & predicted values

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

smaller RSS \rightarrow better model fit

Adjusted R squared

Why adjust?

- penalize adding predictors that don't improve model sufficiently

$$\text{adjusted } R^2 = 1 - \left[\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \times \frac{\frac{n-1}{n-k-1}}{\downarrow \text{new term}} \right]$$

interpretation:

- can decrease when irrelevant predictors added
- better for comparing models with diff # predictors
- preferred over normal R^2 when doing model selection

What is a good R^2 value?

- context/domain specific \rightarrow depends!

high R^2 doesn't mean model is good unless residuals are well behaved + assumptions hold

MLR ASSUMPTIONS

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad i=1, \dots, n$$

for model $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$

ASSUME

L
I
N
E

- ① linearity
- ② independence of errors
- ③ normality of errors
- ④ equal variance of errors

Linearity assume relationship b/w predictors & response is linear

plot residuals vs each predictor (or vs fitted values)

expect: no pattern, randomly scattered

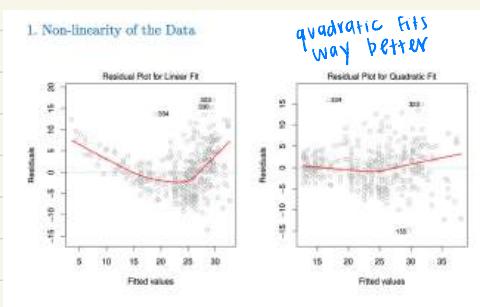
→ curved patterns indicate non-linearity and that a predictor's effect hasn't been fully captured

if violated

→ consider transformation in predictor var

(such as $\log(x)$, quadratic terms,

consider interpretation of transformed vars)



Independence of Errors : assume residuals are indep. of one another

how to check?

- plot residuals vs index of observation (or vs fitted values) → should look random

think about study design

if violated → consider a different model

Normality of Errors : assume residuals are normally distributed

how to check?

- use Q-Q plot of standardized residuals

→ expect points should fall approx. along the 45° line
deviations suggest skewness or heavy tails

* least critical assumption

* doesn't affect predictions, affects inference (p values, confidence intervals)

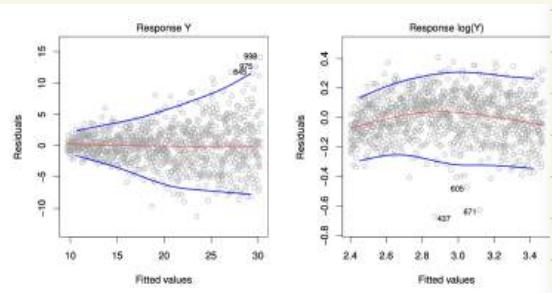
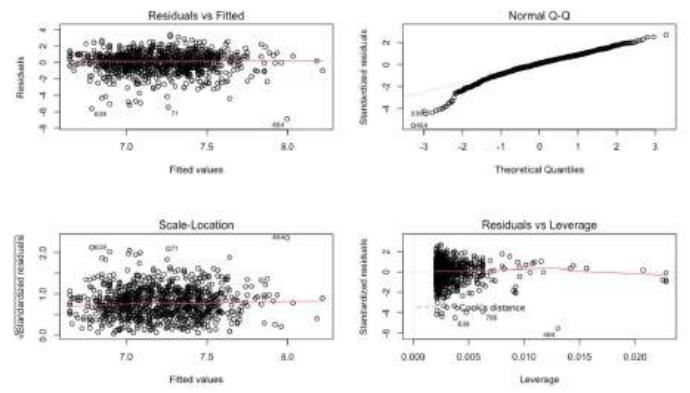
Equal variance/Homoscedasticity: Assume spread of residuals is constant across all levels of predictor
how to check?

- plot residuals vs fitted values
- plot residuals vs index of observation

expect: residuals should be equally spread around zero
→ no fan/funnel shape

if violated

- transform response variable (\log)
- use weighted least squares estimation



Influential Points

Outlier = data pt whose value doesn't follow general trend of rest of data for given var
- investigate as potential data errors / implausible values

Influential points

- individual obs can have large impact on model (estimates, SE, R², RSE)
i.e. a point has large effect on regression

Idea: influential points not always visually obvious from EDA \Rightarrow diagnose with quantitative tools
(leverage scores, Cook's distance)

Leverage = measure how far values of independent vars for i-th obs are from values of other obs

- points with extreme predictor/covariate/feature values are high leverage points

leverage score = i-th diagonal element of

$$H = X [X^T X]^{-1} X^T$$

predictors

matrix to project
observed responses y
into predicted \hat{y}

$0 \leq h_{ii} \leq 1$ and $\sum_{j=1}^n h_{jj} = p + 1$

sum of leverages across
all obs

high leverage points
are obs whose x values
are far from mean of other
 x values

What to do if leverage is high?

- make sure they don't result from data entry errors
- look at impact of those points on estimates

\hookrightarrow high leverage pt doesn't necessarily mean it will have a large effect on regression

influence depends on both x (predictors) and y (response)

example: a point far away in x but close to regression \Rightarrow ~~influential~~ but if regression line \Rightarrow influential

Cook's Distance

- quantify influence of i-th observation
 - idea: remove given obs from dataset
- \hookrightarrow measures how far on avg. predicted values move after removing obs

$$\hat{y}_{j(i)} : \text{predicted value after excluding } i\text{-th obs}$$
$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2}{S_e^2 (p+1)}$$

$\overbrace{S_e^2}^{\text{mean sq error}} \underbrace{(p+1)}_{\text{\# params}}$

$D_i > 1$: strong indication that obs is influential value

$D_i > 0.5$: worthy of attention in smaller datasets

If D_i (Cook's distance) is high, fit model with & without that obs & compare results

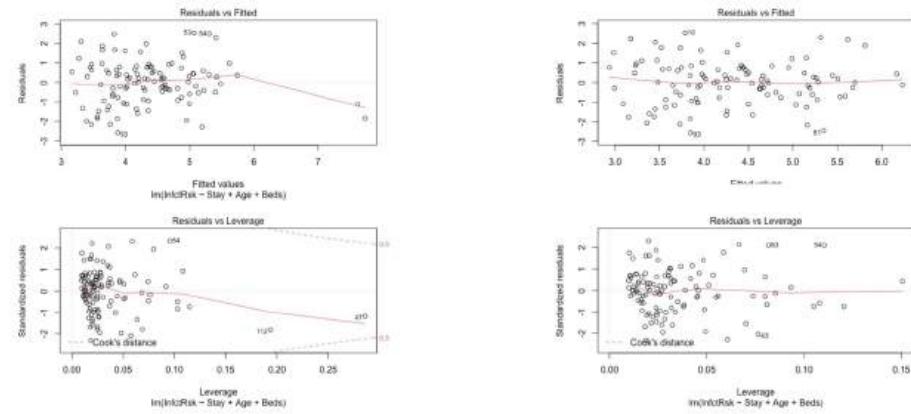
\hookrightarrow are coeffs very different?

\hookrightarrow does model fit change substantially?

*no hard threshold for Cook's distance — consider relative to other pts & look at residual/leverage plots *

Diagnostic plots in R

plot(model, which=5) gives residuals vs leverage plot



what to do with outliers/ influential observations?

- make sure obs is not a data entry error \rightarrow if so, change or exclude
- report results with & without influential obs
 - with large sample size, may not see a diff.
- be wary of numeric cutoffs

multicollinearity

what is multicollinearity?

- when two or more predictor vars in regression model are highly linearly related if two vars are perfectly collinear (perfect linear relationship b/w them) the model cannot compute unique regression coeffs

↳ example

Coefficients: (1 not defined because of singularities)				
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.04536	0.11495	0.395	0.694
X1	4.07300	0.11732	34.716	<2e-16 ***
X2	NA	NA	NA	NA

even near-perfect correlation can be problematic

- inflated standard errors
- unstable estimates
- misleading significance tests

consequences of multicollinearity

- coeff. estimates have large standard errors → low precision
- coeff signs/values become counterintuitive (i.e. negative when we expect positive)
- may fail to identify significant predictors b/c noise drowns signal
- makes interpretation difficult : hard to isolate individual effect of each predictor

when does it occur?

- high correlation b/wn predictors $|P| > 0.9$
- Pearson corr. coeff.
- a prediction is a linear (or nearly) combination of others

How to detect?

during EDA

- is one var derived from another?
- do some vars always increase/decrease together

use correlation matrix to look for predictor pairs with high correlation

use VIF

↳ Variance Inflation Factor (VIF)

$$VIF_j = \frac{1}{1 - R^2_{x_j | X_{-j}}}$$

R squared from regressing x_j on all other predictors

always ≥ 1

quantifies how much the variance of a coeff. is inflated due to multicollinearity

because $R^2 \in [0, 1]$ so $(1 - R^2) \in [0, 1]$

interpretation

\rightarrow not correlated

b/wn 1-5 \rightarrow moderately corr.

$> 5 \rightarrow$ highly correlated

$> 10 \rightarrow$ HIGHLY correlated & want to do something about it

What to do if you detect multicollinearity?

- remove one of the collinear predictors
- scale variables (helps interpretability, but doesn't fix correlation issue)
- in large samples, can be unimportant
 \rightarrow do NOTHING