

## Linear Models with R

### Review

→ (brief) review of vectors/matrices, Probability, and Statistics

### Preface

- objective is to learn what methods are available, and more importantly, when they should be applied

### Ch 1: Introduction

#### 1.1: Before You Start

##### Statistical Approach to a Scientific Problem

- Ask a scientific question and formulate a statistical problem
- collect data
- initial/exploratory analysis
- Answer the question (inferential statistics)

##### Problem Formulation

- understand the physical background and the objective
- make sure you know what is wanted
- put the problem into statistical terms

##### Collect Data

- determine the sampling population
- Sampling
  - + is this a random sample? experiment or observational study?
  - + is there missing data? how is it handled?
  - + is there measurement error?
- if possible, always perform designed experiments

##### 1.2: Initial (Exploratory) Data Analysis

- organize data
- display data graphically
- summarize data
- Be alert for unexpected (e.g. scaling, missingness, etc.)

##### 1.3: When to Use Regression Analysis

##### Inferential Statistics

- Estimate parameters • make predictions • test hypotheses
- What did we learn? what is still uncertain? what may have gone wrong?

##### Regression Analysis

- Build a model to "explain" the relationship between a single variable Y and other variables  $X_1, \dots, X_p$ .
- where  $\begin{cases} Y: \text{response (dependent) var, output} \\ X: \text{predictor (independent) var, input, covariates} \end{cases}$
- tp1: SLR tp2: MLR

Remark: Y: unidimensional → univariate regression

##### Goals of Regression Analysis

- describe data
- make predictions
- assessment of (association) effect of predictor(s) on response

Remark: Association does not imply causation

\* Regression analysis (alone) cannot establish causation

##### Types of Variables

- Qualitative (categorical); can't say one is bigger than another
- Quantitative (numerical)
  - discrete counts → continuous measures

##### What we will cover

- $X$ : continuous, discrete, categorical
  - $Y$ : continuous → linear reg.
  - $Y$ : binary → logistic reg.
  - $Y$ : discrete (counts) → Poisson reg.
- Generalized linear regression  
models (GLMs)  
?

uses a linear combination of predictors to explain a dependent var.

random component

linear (systematic) component

link function

##### Emphasis of the course

- practice using generalized linear regression models
- Learn what methods are available and their limitations
- Many examples w/ less mathematical theory
- More intuition w/ less derivation of formulas
- will still learn mathematical foundations behind practical tools

## Ch 2: Estimation

Textbook	Lecture content
2.1: Linear Model	Regression Analysis
2.2: Matrix Representation	Linear Regression Analysis
2.3: Estimating $\beta$	Simple Linear Regression (SLR)
2.4: Least Squares Estimation	Multiple Linear Regression (MLR)
2.5: Example - calculating $\beta$	Statistical Properties of $\hat{\beta}$ (and $\hat{\sigma}^2$ )
2.6: Gauss-Markov Theorem	Goodness of Fit
2.7: Goodness of Fit	Gauss-Markov Theorem ( $\hat{\beta}$ ; BLUE)
2.8: Example (R - done Separately)	

## Regression Analysis

y: response · output

$x = (x_1, \dots, x_p)$  - predictors · input

Goal: model the relationship between  $y$  and  $x_1, \dots, x_p$

Note: "regression to the mean" phenomena refers to observing extreme events tending to be followed by events closer to the avg.

General form: Let  $X = (x_1, \dots, x_p)^T$ ,  $y = (y_1, \dots, y_n)^T$

then  $y = f(X) + \epsilon$

where  $f(\cdot)$ : underlying truth,  $\epsilon$ : (statistical) error

· there is an unknown mean function  $E(Y|X) = f(X)$

+ we want to estimate  $f(X)$  w/ a statistical model (e.g. regression, NN, Lk)

· Error term  $\epsilon \sim Y - E(Y|X)$

+ measures variability w.r.t. the mean

+ assume:  $E(\epsilon|X) = 0$

· Usually we are given a set of data

+ n: sample size (10 obs)

+ denoted  $(x_{11}, \dots, x_{1p}, y_1), \dots, (x_{n1}, \dots, x_{np}, y_n)$

· can represent in matrix form, i.e.

design matrix  $C(X) = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix}$

response vector ( $y$ ):  $(y_1, \dots, y_n)^T$

## Linear Regression Analysis

· There is no way to estimate  $f(\cdot)$  directly given a finite number ( $n < \infty$ ) of samples

+ functions represent infinite-dimensional space

· We put some restrictions/constraints on  $f(\cdot)$ , i.e.

$f(X) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$  for  $n > p$

where  $\beta_j$ : unknown (shape) parameters,  $\beta_0$ : intercept

$\equiv$  regression coefficients

· estimation of  $f(\cdot)$  is reduced to estimation of  $(p+1)$   $\beta_j$ 's

Additional Ex

(1) GLM s.t.  $E(Y|X) = g(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)$

$\uparrow$  known (often nonlinear) function

e.g. logistic function for logistic regression

(2) NN

+ often used  $n > p$

+ requires additional regularization

(what does "linear" mean?)

· Linear model is linear in parameters, not (necessarily) linear in predictors

Prop. Formally, a function  $g$  is linear in  $\beta$  if

$$g(a \cdot \beta + b \cdot \beta^*) = a \cdot g(\beta) + b \cdot g(\beta^*)$$

where  $a, b \in \mathbb{R}$ ,  $\beta, \beta^* \in \mathbb{R}^p$

Corollary: Equivalently, underlying statistics has the following prop's:

· Additivity:  $f(x+y) = f(x) + f(y)$

· Homogeneity of degree 1:  $f(cx) = c f(x) \forall c$ .

Ex:  $f(x) = \beta_0 + \beta_1 x_1 + \beta_2 \ln(x_2) + \beta_3 x_3$  is a linear model

Non-Ex:  $f(x) = \beta_0 + \beta_1 x_1 x_2$  is not

## Transformations

$f(x) = \beta_0 x_1^{\beta_1}$  is not a linear model.

However, notice that

$$\ln(f(x)) = \ln(\beta_0) + \beta_1 \ln(x_1)$$

Hence, if we let  $\tilde{x}_1 = \ln(x_1)$ ,  $\tilde{x}_2 = \ln(\beta_0) + \beta_1 x_1$ , we have

$$f(x) = \beta_0 e^{\beta_1 \ln(x_1)} = \beta_0 e^{\tilde{x}_1}$$

which is a linear model!

## Implications

· Linear models are less restrictive than you might think

· they can be made very flexible by transformation of the response & predictors

Def (linear transformation):  $f(x) = \beta_0 + \beta_1 g_1(x) + \dots + \beta_p g_p(x)$

where  $g_1, \dots, g_p$  are known functions

· Linear models are not necessarily straight lines, e.g.

$$y = \alpha x^2 + bx + c$$

### Simple Linear Regression

Let  $p = \text{number of predictors}$

$$\text{Model: } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon \quad (1)$$

$$\text{s.t. } E(Y|X) = \beta_0 + \beta_1 X_1$$

Assumptions: $E(\epsilon X) \sim N(0, \sigma^2)$	1) $E(E(X)) = 0 \Rightarrow E(Y X) = \beta_0 + \beta_1 X_1$ + zero mean + linearity
	2) $V(E(X)) = \sigma^2 \Rightarrow V(Y X) = \sigma^2$ + constant variance + $V(\beta_0 + \beta_1 X_1 + \epsilon) = \sigma^2$ by var. of const.
	3) Given $X_i$ 's, $\epsilon_i$ 's are iid. (or uncorrelated) $\Rightarrow$ normality

**Remark:** For violations of assumptions, namely (3), there are modifications to model (1) we can make that are discussed later

Let $E = (E_1, \dots, E_n)^T$ , $X = (X_1, \dots, X_n)^T$
Then, equivalently assume: (1) $E(E X) = \begin{pmatrix} E(E_1 X) \\ \vdots \\ E(E_n X) \end{pmatrix} = 0$
(2) $\text{Cov}(E X) = \begin{pmatrix} V(E_1 X) & & & \\ Cov(E_1, E_2 X) & \ddots & & \\ & \ddots & \ddots & \\ & & Cov(E_{n-1}, E_n X) & \\ & & & V(E_n X) \end{pmatrix}_{n \times n} = \begin{pmatrix} \sigma^2 & & & \\ & \ddots & & \\ & & \sigma^2 & \\ & & & \sigma^2 \end{pmatrix} = \sigma^2 I_n$
(3) Stronger assumption assumes $E_i$ 's are independent
(4) Further assume $E_i$ 's are normally distributed

**Goal:** Given  $(x_i, y_i), i \in \{1, \dots, n\}$ , estimate  $\beta_0, \beta_1$ , and  $\sigma^2$

+ we do this by minimizing errors

$$\rightarrow \text{One criterion is least squares: } \min_{\beta_0, \beta_1} \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (2)$$

$$= \sum_{i=1}^n (y_i - E(Y|X_i))^2$$

$$\text{Remark: observe the special case s.t.: } \begin{array}{c} Y \\ \downarrow \text{---} \\ E(Y|X) = \beta_0 \Rightarrow \hat{\beta}_0 = \arg \min_{\beta_0} \sum_{i=1}^n (y_i - \beta_0)^2 \\ \Rightarrow \hat{\beta}_0 = \bar{y} \end{array}$$

Next, we differentiate (2) w.r.t.  $\beta_0, \beta_1$ , set to zero, and solve, i.e.

$$\frac{\partial}{\partial \beta_0} = (-2) \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\Rightarrow E(Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\Rightarrow \hat{\beta}_0 + \frac{1}{n} \sum_{i=1}^n E(Y_i - \hat{\beta}_1 x_i) = 0$$

$$\Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

and

$$\frac{\partial}{\partial \beta_1} = (-2) \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\Rightarrow E(x_i Y_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2) = 0$$

$$\Rightarrow E(x_i Y_i) - \hat{\beta}_0 E(x_i) - \hat{\beta}_1 E(x_i^2) = 0$$

$$\Rightarrow E(x_i Y_i) - (\bar{y} - \hat{\beta}_1 \bar{x}) \bar{x} - \hat{\beta}_1 E(x_i^2) = 0, \text{ substitution}$$

$$\Rightarrow E(x_i Y_i) - n \bar{y} \bar{x} = \hat{\beta}_1 E(x_i^2) - \hat{\beta}_1 n \bar{x}^2$$

$$\Rightarrow \hat{\beta}_1 = \frac{E(x_i Y_i) - n \bar{y} \bar{x}}{E(x_i^2) - n \bar{x}^2} \quad (3)$$

where

a) is equivalent to  $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ :

$$\begin{aligned} \sum (x_i - \bar{x})(y_i - \bar{y}) &= \sum (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y}) \\ &= \sum x_i y_i - \bar{y} \sum x_i - \bar{x} \sum y_i + \bar{x} \bar{y} \\ &= \sum x_i y_i - \bar{y} (n \bar{x}) - \bar{x} (n \bar{y}) + n \bar{x} \bar{y} \\ &= \sum x_i y_i - n \bar{x} \bar{y} \end{aligned}$$

and b) is equivalent to  $\sum_{i=1}^n (x_i - \bar{x})^2$ :

$$\begin{aligned} \sum (x_i - \bar{x})^2 &= \sum (x_i^2 - 2 \bar{x} x_i + \bar{x}^2) \\ &= \sum x_i^2 - 2 \bar{x} \sum x_i + n \bar{x}^2 \\ &= \sum x_i^2 - 2 \bar{x} \bar{x} + n \bar{x}^2 \\ &= \sum x_i^2 - 2 n \bar{x}^2 + n \bar{x}^2 \\ &= \sum x_i^2 - n \bar{x}^2 \end{aligned}$$

This implies that:

$$\left\{ \begin{array}{l} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{array} \right.$$

Remarks: - not notation used for estimates

$$\cdot \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\Leftrightarrow \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

$\Rightarrow (\bar{x}, \bar{y})$  lies on the estimated regression line.

$$\cdot \hat{\beta}_1 = \frac{\text{cov}(X, Y)}{V(X)} \quad \left| \begin{array}{l} \text{P: cov}(X, Y) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ V(X) = \sum_{i=1}^n E(X_i - \bar{X})^2 = S_x^2 \end{array} \right. \Rightarrow \frac{\text{cov}(X, Y)}{S_x^2} = \frac{\text{cov}(X, Y)}{V(X)} \Rightarrow \hat{\beta}_1 = \frac{\text{cov}(X, Y)}{S_x^2}$$

$$\cdot \hat{\beta}_0 = \frac{\text{cov}(X, Y)}{V(X)} = \frac{\text{cov}(X, Y)}{S_x^2} \frac{S_y}{S_x} \quad \left| \begin{array}{l} \text{P: cov}(X, Y) = \frac{\text{cov}(X, Y)}{S_x S_y} \Rightarrow \text{cov}(X, Y) = \text{cov}(X, Y) S_x S_y \\ \text{so, } \frac{\text{cov}(X, Y)}{V(X)} = \frac{\text{cov}(X, Y) S_x S_y}{S_x^2} = \text{cov}(X, Y) \frac{S_y}{S_x} \end{array} \right.$$

### Geometric Interpretation

Formally, let  $r = \text{cor}(X, Y)$ ,  $S_y = \text{SD}(Y)$ ,  $S_x = \text{SD}(X)$ , we can rewrite

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

$= \bar{y} + \hat{\beta}_1 (\bar{x} + x - \bar{x})$  via substitution

$$\Leftrightarrow y - \bar{y} = \hat{\beta}_1 (\bar{x} + x - \bar{x})$$

$$= r \frac{S_y}{S_x} (\bar{x} + x - \bar{x})$$

$$\Leftrightarrow \frac{y - \bar{y}}{S_y} = r \frac{(\bar{x} + x - \bar{x})}{S_x}$$

or, if  $x$  and  $y$  are standardized first (mean 0 SD 1), simply

$$y = r \frac{x}{S_x} \quad \left| \begin{array}{l} \text{P: } \hat{\beta}_1 = \frac{\text{cov}(X, Y)}{V(X)} = \frac{\text{cov}(X, Y) S_x S_y}{V(X) S_y} = \text{cor}(X, Y) \text{ when } x, y \text{ are standardized.} \end{array} \right.$$

recall: standardizing; for  $i=1, \dots, n$

$$\begin{aligned} x_i &\rightarrow \frac{x_i - \bar{x}}{S_x} = x_i^* \\ y_i &\rightarrow \frac{y_i - \bar{y}}{S_y} = y_i^* \end{aligned} \quad \left| \begin{array}{l} \text{3: } x_i^* \text{ and } y_i^*, i=1, \dots, n \\ \text{4: sample means 0 and SD 1} \end{array} \right.$$

### Two regression lines

Suppose  $x$  and  $y$  have both been standardized

- regress  $y$  on  $x$ :  $y \sim x$  when  $|r| \leq 0.5 - 1 \leq r \leq 1$   
 - regress  $x$  on  $y$ :  $x \sim y$

regression effect: predictions always "regress" towards the mean

@: using  $x \sim y$  and not  $x = \frac{y}{r}$   
 + two regressions minimize different error directions

### Multiple Linear Regression

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

Remarks: -  $x_i \in \mathbb{R}^{(p+1)}$  =  $(1, x_{i1}, \dots, x_{ip})^T$

-  $\beta \in \mathbb{R}^{(p+1) \times 1} = (\beta_0, \beta_1, \dots, \beta_p)^T$

-  $E(Y|X_1, X_2) = \beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi}$

( $\Rightarrow E(Y|X_1, X_2) \neq x_1^* \beta$ )

-  $p$  predictors

( $\Rightarrow p$  parameters and  $(1) \sigma^2 = \text{var}(\varepsilon_i | x_i)$  parameter)

Assumptions: 1) (linear) mean function  $E(Y|X_1, X_2) = x_i^* \beta$  ( $\Rightarrow E(E(Y|X_1, X_2)) = 0$ )  
 $\varepsilon_i \sim N(0, \sigma^2)$  2) var function  $\text{var}(Y|X_1, X_2) = \sigma^2$  ( $\Rightarrow \text{var}(\varepsilon_i | x_i) = \sigma^2$  for  $i=1, \dots, n$ )  
 3)  $E(\varepsilon_i | x_i)$  independent or uncorrelated ( $\text{cov}(E(\varepsilon_i | x_i), E(\varepsilon_j | x_j)) = 0$  for  $i \neq j$ )  
 $\Leftrightarrow$  4)  $E(E(\varepsilon_i | x_i)) = 0$  for  $E_{\varepsilon|x} = (E_1, \dots, E_n)^T$   
 5)  $\text{cov}(E_{\varepsilon|x}) = \sigma^2 I_n$

### Matrix Notation

Let  $y = (y_1, \dots, y_n)^T$ ,  $X \in \mathbb{R}^{n \times (p+1)}$   $\left( \begin{array}{c} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{array} \right) = \left( \begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right)$ ,  $\beta \in \mathbb{R}^{(p+1)}$   $\left( \begin{array}{c} \beta_0 \\ \vdots \\ \beta_p \end{array} \right)$ ,  $\varepsilon \in \mathbb{R}^n$   $\left( \begin{array}{c} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{array} \right)$

Remarks: -  $X \in \mathbb{R}^{n \times (p+1)}$  column

. 1<sup>st</sup> row of  $X$  is  $x_1^T$  where  $x_1^T = \left( \begin{array}{c} x_{11} \\ \vdots \\ x_{1p} \end{array} \right)$ ,  $i=1, \dots, n$

thus, more compactly

$$y = X \beta + \varepsilon \quad \Leftrightarrow \left( \begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right) = \left( \begin{array}{c} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{array} \right) \left( \begin{array}{c} \beta_0 \\ \vdots \\ \beta_p \end{array} \right) + \left( \begin{array}{c} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{array} \right)$$

### Estimating $\beta$

- observe  $y$  and  $X$  (n samples)

- want to minimize errors

Least Squares Criteria: min  $\sum_i \varepsilon_i^2 = \varepsilon^T \varepsilon$ ,  $\varepsilon = y - X\beta$   
 $\Rightarrow \min \sum_i (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$   
 $\Rightarrow (y - X\beta)^T (y - X\beta)$ ,  $\varepsilon^T \varepsilon$

$$= y^T y - 2y^T X\beta + \beta^T X^T X\beta$$

where (cross product term):  $-2y^T X\beta = -2(y^T x_1)\beta_1 - \dots - 2(y^T x_p)\beta_p$

$$= -2y^T x_1 \beta_1 - \dots - 2y^T x_p \beta_p$$

$$= -2y^T x \beta$$

Differentiating (7) wrt  $\beta$  and setting to 0:

$$\begin{aligned} 0 &= \eta_p(C(X^T X)\beta + C^T X^T y) \\ &\leftarrow -2(C^T X)^T + 2C(X^T X)\beta, \text{ by } (1), (2) \text{ resp.} \\ 0 &= -2X^T y + 2C(X^T X)\beta \\ \Rightarrow X^T y &= C(X^T X)\beta \end{aligned}$$

**Normal equation:**  $X^T y = C(X^T X)\beta$

(\*)

**ordinary least squares (OLS) estimator**

$$\hat{\beta} = C(X^T X)^{-1} X^T y$$

- assumes  $X^T X$  invertible

-  $X_{\text{ncpt}}$  full column rank

↳  $X^T X$  invertible

notes: (if  $X^T X$  not invertible)

**matrix-calculus formulas**

Let  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times n}$ ,  $C \in \mathbb{R}^{n \times n}$  (symmetric)

recall:  $\text{tr}(A) = a_{11} + a_{22} + \dots + a_{nn}$

$$1) \nabla_B \text{tr}(AB) = A^T$$

$$2) \nabla_B \text{tr}(C^T B) = 2CB$$

(\*) when  $m=1 \Leftrightarrow A \in \mathbb{R}^n$ ,  $B \in \mathbb{R}^n$  vectors

Note that  $\text{tr}(AB) = A^T B$

$$\text{tr}(B^T B) = B^T B$$

Now,

$$1) \nabla_B (AB) = A^T$$

$$2) \nabla_B (B^T CB) = 2CB$$

**Remark:**  $\|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2} = \sqrt{\text{tr}(A^T A)} = \sqrt{\text{tr}(A^2)} = \sqrt{\sum_{i,j} |a_{ij}|^2}$ ,

i.e. the Frobenius norm: standard Euclidean norm of a matrix treated as a long vector.

**Remarks:** - we discuss alternative methods for when  $n < p+1$  (over)

- we also discuss (multi)collinearity, multicollinearity, isolation later

### Error model

$$\begin{aligned} \text{FORM: } & \text{value } E(Y_i | X_i) = \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip} \\ \text{MODEL: } & E(Y_i | X_i) = f(X_i) = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip} \end{aligned}$$

$$\text{RESIDUALS: } \hat{e}_i = e_i = y_i - \hat{y}_i = y_i - f(x_i)$$

**Residual vector:**  $\hat{e} = (\hat{e}_1, \dots, \hat{e}_n)^T$ ;  $\hat{e}_{\text{tot}} = Y_{\text{tot}} - \hat{Y}_{\text{tot}} = Y - X\hat{\beta}$

**Properties:** (1) sample mean of residuals is zero

$$\Rightarrow \bar{e} = \frac{1}{n} \sum_i \hat{e}_i = 0$$

$$\Rightarrow X^T \hat{e} = X^T Y - X^T X\hat{\beta} = 0 \quad (2) \quad X^T (Y - X\hat{\beta}) = 0 \quad (3) \quad X^T \hat{e} = 0$$

First row vector of  $X^T$  is  $(1, \dots, 1)_{1 \times n}$

$$\Rightarrow (1, \dots, 1) \hat{e}_{\text{tot}} = 0 \Rightarrow \frac{1}{n} \hat{e}_i = 0 \Rightarrow \bar{e} = 0.$$

**Remarks:** Normal eq. further implies that residual and predictor are uncorrelated to each other, i.e.  $X^T \hat{e} = 0$ .

$$\text{RESIDUAL SUM OF SQUARES (RSS): } \sum_{i=1}^n \hat{e}_i^2$$

### Hat Matrix

$$\hat{y} = X\hat{\beta} = X(C(X^T X)^{-1} X^T) y = Hy \text{ s.t.}$$

$$\boxed{\text{HAT MATRIX } H = X(C(X^T X)^{-1} X^T)}$$

Fitted values:  $\hat{y} = Hy$

$$\text{residuals: } \hat{e} = y - \hat{y} = y - (Hy) = (I - H)y$$

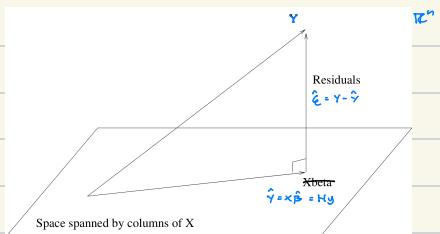
$\Rightarrow H$  is a projection matrix

### Projection matrix

- [Def:]  $H$  is a projection matrix if
  - ↳  $H^T = H$  (symmetric)
  - ↳  $HH^T = H$  (idempotent)

note: able to show that  $X(C(X^T X)^{-1} X^T)$  satisfies (i) and (ii)

### Vector Space Interpretation



$R^n$

\* the projection matrix  $H$  projects  $y$  onto the column space of  $X_{\text{ncpt}}$ , which leads to the vector space interpretation of least squares estimate - subspace spanned by column vectors of  $X_{\text{ncpt}}$

(\*) Let  $C(X) :=$  column space of  $X$ .

$$\text{then } P_{C(X)}(y) = Hy = \hat{y} = X\hat{\beta} = X(C(X^T X)^{-1} X^T) y, \text{ i.e.}$$

the hat matrix ( $H$ ) is the projection operator  $P_{C(X)}$

$$\text{s.t. } H: y \rightarrow \hat{y}$$

\*  $\min_{\beta} \|y - X\beta\|^2$  minimizes the Euclidean dist. b/w  $y$  and the linear space spanned by columns of  $X$

$$\hat{e} \perp \text{linear space} \Leftrightarrow X^T \hat{e} = 0$$

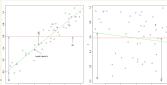
$$\hat{e} \perp Y \Leftrightarrow X^T \hat{e} = 0$$

$$\Leftrightarrow X^T X \hat{e} = 0$$

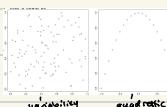
$$\Leftrightarrow X^T \hat{e} = 0.$$



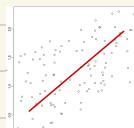
- $R^2$  closer to 1 indicates good fit



- $R^2 \approx 0$  for different reasons



- Small  $R^2$  does not necessarily mean that  $y$  and  $X$  are not linearly related; can have slight trend w/ high variance



$$R^2 \approx 0.2$$

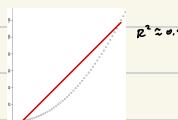
- In SLR:  $R^2 = r^2$

$$R^2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \text{ and } \hat{y}_i = \bar{y} + \hat{p}_i(x_i - \bar{x}).$$

$$\text{Now, } SSe_{reg} = \sum (\hat{y}_i - \bar{y})^2 = \sum (\hat{p}_i(x_i - \bar{x}))^2 = \hat{p}_i^2 \sum (x_i - \bar{x})^2 \\ = \left( \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right)^2 \sum (x_i - \bar{x})^2 = \frac{[\sum (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum (x_i - \bar{x})^2}.$$

- $R^2 \approx 1$  does not mean the linear model is correct

$$\text{Then, } R^2 = \frac{[\sum (x_i - \bar{x})(y_i - \bar{y})]^2}{[\sum (x_i - \bar{x})^2][\sum (y_i - \bar{y})^2]} = \left[ \frac{\text{cov}(x, y)}{s_{\text{err}}^2} \right]^2 = r^2. \quad \square$$



### Adjusted $R^2$

adjust for the number of predictors in the model

$$\text{adj-}R^2 = 1 - \frac{RSS/(n-p-1)}{TSS/(n-1)} \\ = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

REMARKS: - Penalizing (correction) factor for additional predictors

→ helps w/ model selection/comparison value  $R^2$

(or stay the same)

- $R^2$  will always increase when adding more predictors
- while adj- $R^2$  will not necessarily increase (probabilistic chance - but high likelihood)

### Gauss-Markov Theorem

- Q: why use the least-squares estimate  $\hat{\beta}$ ?

DEF:  $\hat{\beta}$  estimate  $\beta$ , we say  $\hat{\beta}$  is unbiased if  $E(\hat{\beta}) = \beta$

DEF: we say  $(\hat{\beta})$  is linear if it is a linear function in  $y$ :  
e.g. OLS  $\hat{\beta} = (X'X)^{-1}X'y$  is linear

THEM: Suppose  $y = XB + \epsilon$ ,  $X$  full column rank  
(\*)  $E(\epsilon) = 0$ ,  $\text{Var}(\epsilon) = \sigma^2 I$ .

consider  $\gamma = c^T \beta$ , i.e. any linear combination of unknown  $\beta_0, \beta_1, \dots, \beta_p$  and  $c$  is a known, pre-specified vector.

Then, among all unbiased linear estimates of  $\gamma$ ,  $\hat{\gamma} = c^T \hat{\beta}$  has the minimum variance and is unique

REMARKS: Let  $C = (1, x_1, \dots, x_p)$ , then  $\gamma = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

OLS solution is the Best Linear Unbiased Estimate (BLUE)

$E(\hat{\gamma}) = \gamma \Rightarrow E(C^T \hat{\beta}) = C^T \beta$

$\hat{\beta}$  unbiased and linear estimator  $\hat{\beta}$ ,

$$\text{Var}(C^T \hat{\beta}) \geq \text{Var}(C^T \beta)$$

where  $\hat{\beta}$ : OLS and equality holds, i.e.  $\hat{\beta} = \beta$ .

Situations where estimators other than OLS should be used:

+ correlated errors and/or unequal variance: generalized least squares (ch. 1)

+ long-tailed error distribution: use robust estimators (ch. 1), typically not linear in y

+ collinear predictors: try biased estimators, e.g. ridge regression (ch. 9)

PROOF: Let  $\tilde{\beta}$  be any other linear unbiased estimator for  $\beta$ .

then  $\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta})$  and  $\text{Var}(\tilde{\beta}) = \text{Var}(\hat{\beta}) \Leftrightarrow \tilde{\beta} = \hat{\beta}$ .

Since  $\tilde{\beta}$  is a linear estimator, it can be written as  $\tilde{\beta} = d^T y$  for  $d \in \mathbb{R}^n$ ,  $d \neq 0$ .  
Thus,  $\tilde{\beta} = d^T y = d^T X \beta + d^T \epsilon$ .

Note that  $E(\tilde{\beta}) = d^T \beta$  and since  $\hat{\beta}$  is unbiased,  $d^T \beta = c^T \beta$ .

Since this is true for any  $\beta$ , it must be that  $d^T = c^T$ .

Now let  $\rho = d - c^T (X'X)^{-1}X'$ . Note

$$\rho^T = d^T - c^T (X'X)^{-1}X^T$$

and therefore, multiplying both sides by  $X$ ,

$$\rho^T X = d^T X - c^T (X'X)^{-1}X^T X$$

$$= d^T X - c^T$$

$$= 0$$

Next note  $d^T = c^T (X'X)^{-1}X'$  and therefore

$$\text{Var}(\tilde{\beta}) = \text{Var}(d^T \epsilon)$$

$$= \sigma^2 d^T d$$

$$= \sigma^2 [c^T + c^T (X'X)^{-1}X^T] [c^T (X'X)^{-1}X^T]$$

$$= \sigma^2 [c^T c + c^T (X'X)^{-1}X^T]$$

$$= \sigma^2 c^T c + \text{Var}(\hat{\beta})$$

$$\geq 0 \text{ and } \rho^T \rho = 0 \Rightarrow \rho = 0 \Rightarrow \tilde{\beta} = \hat{\beta}.$$

## What can go wrong?

- $X^T X$  could be singular (chappas if predictors are linearly dependent or if  $p=n$ ).
- The following assumptions could be violated:
 
$$\begin{cases} \text{Var}(\epsilon) = \sigma^2 I & (\text{constant var}) \\ E(\epsilon) = 0 & (\text{uncorrelated}) \end{cases}$$

Note: OLS best only among linear, unbiased estimates

## Ch. 3: Inference

### Textbook

- 3.1: HT's to compare models
- 3.2: Testing Examples
- 3.3: Permutation Tests (SKIDP)
- 3.4: (Simultaneous) CI's for  $\beta$
- 3.5: CI's for Prediction
- 3.6: Designed Experiments
- 3.7: Observational Data
- 3.8: Practical Difficulties

### Lecture content

- Hypothesis Testing
- $t$  distribution
- (General) F-test and/or t-test
- Testing more than one Predictor
- Confidence Intervals
- Simultaneous Confidence Regions
- CI's for Prediction
- Interpreting Parameter Estimates
- Experimental Design
- Randomization
- Orthogonality

## Inference

Estimates:  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$

- want to draw conclusions about  $\beta_0, \beta_1, \dots, \beta_p$
- two main inference tools (HT and CI)

## Hypothesis Testing

- use probability to decide whether data is consistent w.r.t hypothesis

Null Hypothesis:  $= H_0$  (e.g.  $\beta_0 = 0$ )

Alternative Hypothesis:  $= H_1$  (e.g.  $\beta_0 \neq 0$ )

- Decide whether data is consistent w.r.t  $H_0$ :  $\begin{cases} \text{if not } \rightarrow \text{reject } H_0 \text{ and accept } H_1 \\ \text{otherwise } \rightarrow \text{FTR } H_0 \end{cases}$

## Errors in Hypothesis Testing

		True State		TPE := $P(\text{FTR } H_0   H_0 \text{ true})$ : false positive (error)	TNE := $P(\text{CTFR } H_0   H_1 \text{ true})$ : false negative (error)
		$H_0$ true	$H_1$ true		
Decision	Accept $H_0$	✓	TPE		
	Reject $H_0$	TPE	✓		

(Statistical Power) :=  $1 - TNE = 1 - \text{FNR} = \text{true positive (rate)}$   
 $\text{CNR} = (TPE)^2 = P(\text{reject } H_0 | H_0 \text{ false})$ ; correctly detect effect when one exists

## Procedure

Set significance level:  $\alpha = P(TPE)$

→ typically  $\alpha = 0.05$  (alpha-level), 0.01

Compute the p-value := prob of observed or more extreme departure from  $H_0$  (in favor of  $H_1$ ) when  $H_0$  is true

p-value  $< \alpha \Rightarrow \text{reject } H_0$

## Further Assumption on Errors

In addition to linearity,  $E(\epsilon) = 0$ , we assume a distribution for the errors, i.e.

$$\epsilon \sim N(0, \sigma^2 I_n) \quad (\Leftrightarrow E(\epsilon|I) = 0, \text{covar}(\epsilon) = V(\epsilon|I) = \sigma^2 I_n)$$

## Distribution of $\hat{\beta}$

recall:  $\hat{\beta} = \beta + C(X^T X)^{-1} X^T \epsilon$

If  $\epsilon \sim N(0, \sigma^2 I_n)$ , then

$$\begin{aligned} \hat{\beta} &\sim N_{p+1}(C(\beta), (X^T X)^{-1} \sigma^2 I) \\ \hat{\beta} &\sim N(C\beta, (X^T X)^{-1} \sigma^2 I) \end{aligned}$$

Where, in practice, we use the approximation S.E.

$$\hat{S.E.}(\hat{\beta}_j) := \sqrt{\frac{\sigma^2}{n} \frac{(X^T X)^{-1}_{jj}}{(n-p)}} \quad \Rightarrow \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-(p+1)} = \frac{RSS}{n-(p+1)}$$

Under the normal assumption on the errors:

$$\begin{aligned} \hat{\beta}_j / \hat{S.E.}(\hat{\beta}_j) &\sim N(0, 1) \quad \text{and} \\ \hat{\beta}_j - \beta_j / \hat{S.E.}(\hat{\beta}_j) &\sim t_{n-(p+1)} \end{aligned}$$

$$\text{Proof (Sketch): } \hat{\beta}_j - \beta_j \stackrel{d}{\sim} N(0, 1) \quad \text{and} \quad \hat{S.E.}(\hat{\beta}_j) \stackrel{d}{\sim} \frac{\sum_{i=1}^{n-p} z_i^2, z_i \stackrel{iid}{\sim} N(0, 1)}{n-(p+1)} = \frac{\hat{\sigma}^2}{n-(p+1)} \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\hat{S.E.}(\hat{\beta}_j)} \stackrel{d}{\sim} t_{n-(p+1)}.$$

$$\text{recall: } t_{df} := \frac{N(0, 1)}{\sqrt{\frac{\sigma^2}{df}}} \quad (\Leftrightarrow t_{n-(p+1)} = \frac{N(0, 1)}{\sqrt{\frac{\sigma^2}{n-(p+1)}}} \quad \text{where } \chi^2_{n-(p+1)} \stackrel{d}{\sim} \sum_{i=1}^{n-p} Z_i^2, Z_i \stackrel{iid}{\sim} N(0, 1))$$

## t-distribution

Probability density function (pdf):

$$\begin{aligned} N(0, 1) &\sim \frac{1}{\sqrt{2\pi}} \exp^{-\frac{x^2}{2}} \\ t_n &\sim \Gamma\left(\frac{n+1}{2}\right) \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} / \sqrt{\pi \Gamma\left(\frac{n}{2}\right)} \end{aligned}$$

Remarks: t has a single parameter  $n$  (df)

• symmetric around 0, "less skewed", heavier tails than normal

• as  $n \rightarrow \infty$ ,  $t_n \rightarrow N(0, 1)$

•  $P(|t_{df}| > t) = \int_{-t}^t \frac{1}{\sqrt{\frac{\sigma^2}{df}}} \exp^{-\frac{x^2}{2}} dx$

•  $P(|t_{df}| > t) = \int_{-t}^t \frac{1}{\sqrt{\frac{\sigma^2}{df}}} \exp^{-\frac{x^2}{2}} dx$

### Another (General) Approach: F-test

recall:  $RSS = \sum_i e_i^2$

remains:  $H_0$  model is nested (simpler, special case) to the  $H_A$  model

• Fit a model under  $H_0$  and compute  $RSS_{H_0}$  (e.g.  $\beta_0, \beta_1, \dots = 0$ )

$$\text{e.g. } X \text{ pop: } k_2, H_0: \beta_2 = 0 \Leftrightarrow E(Y|X) = \beta_0 + \beta_1 X_1 + \beta_3 X_3 + \beta_4 X_4$$

• Fit another model under  $H_0 \vee H_A$  and compute  $RSS_{H_0 \vee H_A}$  (e.g. no restriction on  $\beta$ 's)

$$\text{e.g. } H_0: \beta_2 = 0 \Leftrightarrow E(Y|X) = \beta_0 + \beta_1 X_1 + \beta_3 X_3 + \beta_4 X_4$$

• Now compute

$$F = \frac{\left[ \frac{RSS_{H_0} - RSS_{H_0 \vee H_A}}{df_{H_0} - df_{H_0 \vee H_A}} \right]}{\left[ \frac{RSS_{H_0 \vee H_A}}{df_{H_0 \vee H_A}} \right] (\#)}$$

where  $(\#) = \hat{\sigma}^2$  under  $H_0 \vee H_A$  models

If  $H_0$  is true

$$F \sim F_{df_1, df_2} \text{ where } \begin{cases} df_1 = df_{H_0} - df_{H_0 \vee H_A} \\ df_2 = df_{H_0 \vee H_A} \end{cases}$$

compute p-value =  $P(F_{df_1, df_2} > F)$

### Summary

1. Define (larger) full model (FM); thought to be most appropriate for data
2. Define (smaller) reduced model (RM); described by  $H_0$
3. Use F-test to decide whether not to reject smaller model in favor of  $H_0$  (larger, full model)

$H_0: R.M. = M_{H_0}, \beta_i = 0$  for  $i$ : number of parameters in FM,  $R.M.$

$H_A: M_{H_0} \neq M_{H_A}, \beta_i \neq 0$

↳ reject  $H_0$

↳ there is (at least) one significant linear relationship between  $\beta_i$  and  $Y$

### F-distribution

Let  $Z_1, \dots, Z_n \stackrel{iid}{\sim} N(0, 1)$ . Then

$$U = Z_1^2 + \dots + Z_n^2 \sim \chi^2_n = \text{Gamma}\left(\frac{n}{2}; \frac{1}{2}\right).$$

Prob: Suppose  $U \sim \chi^2_n, W \sim \chi^2_m$  are independent. Then

$$\frac{U/n}{W/m} \sim F_{n,m}$$

Remarks:

$$\begin{cases} F_{df_1, df_2} \geq 0 \\ F_{df_1, df_2} \sim F_{df_1, df_2} \end{cases} \quad P(\text{when not } H_0, U = z, z) \Rightarrow \frac{U/n}{W/m} = \frac{z^2}{W/m} = \left(\frac{z}{\sqrt{W/m}}\right)^2 \sim t_{m-1}^2,$$

(ANOVA ( $M_{H_0}, M_{H_A}$ )  $\rightarrow$  Res.DF RSS DF Sum of Square F Pr(>F))

$$\begin{array}{ll} \text{Res.DF} & \text{RSS} \\ \text{df}_{H_0} & \text{RSS}_{H_0} \\ \text{df}_{H_A} & \text{RSS}_{H_A} \end{array} \quad \text{df}_{H_0} = \text{df}_{H_A} \quad \text{RSS}_{H_A} = \text{RSS}_{H_0} \quad \text{F-stat: } \text{Pr}(>F)$$

where  $H_0: M_{H_0} = M_{H_A}$

### F-test and t-test

Prob: F-test and two-sided t-test are equivalent for testing a single predictor

Prob: Consider SCM S.I.  $E(Y|X) = \beta_0 + \beta_1 X_1, X \in \mathbb{R}$

Show that F-test and t-test are equivalent for testing

$H_0: \beta_1 = 0$  vs.  $H_A: \beta_1 \neq 0$

i.e. prove  $\Rightarrow$  F-test = t-test

$$\Leftrightarrow \frac{\left[ \frac{RSS_{H_0} - RSS_{H_0 \vee H_A}}{df_{H_0} - df_{H_0 \vee H_A}} \right]}{\left[ \frac{RSS_{H_0 \vee H_A}}{df_{H_0 \vee H_A}} \right]} = \left[ \frac{\bar{t}_i^2}{SE(\beta_i)} \right]^2$$

$$\Rightarrow F_{1, n-2} \sim (t_{n-2})^2 \quad \square$$

### Testing more than one Predictor

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

$H_A: \text{not } H_0$ , i.e. at least one  $\beta_i \neq 0$  for  $i=1, \dots, p$

↳ Testing a Pair

$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \Leftrightarrow E(Y|X) = \beta_0$

$H_A: \text{not } H_0$ ; i.e. at least one  $\beta_i \neq 0$  for  $i=1, \dots, p$

↳ Overall F-test

(as seen in overall R Summary table)

Suppose we want to test the equality of two effects:

$H_0: \beta_1 = \beta_2 \Leftrightarrow \beta_1 X_1 + \beta_2 X_2 = \beta_0 (X_1 + X_2)$ ; (so  $L = \{Y = \beta_0 (X_1 + X_2) + \dots\}$ )

$H_A: \beta_1 \neq \beta_2$

↳ testing a subspace

Suppose we want to test an effect value, i.e.

$$H_0: \beta_0 = 0 \Rightarrow \text{Im}(y - \beta_0 - \dots - \text{offset}(\text{cov}(X)))$$

$$H_1: \beta_0 \neq 0$$

(F-RST approach)

$\Leftrightarrow$  same  $H_0$  and  $H_1$

$$\text{w/ } t_{\text{stat}} \sim \frac{\hat{\beta}_0 - (\beta_0 = 0)}{\text{SE}(\hat{\beta}_0)}$$

testing a subspace

### confidence intervals

H.T comments:

- only answers yes/no questions
- dependence on sample space  
+ larger sample space  $\rightarrow$  more likely to reject  $H_0$
- notion of statistical significance vs. practical significance

e.g. assuming standardized data

$$\hat{\beta} = 0.001, \text{ testing } H_0: \beta_0 = 0 \quad \text{vs. } p < 0.001 \text{ statistically significant}$$

$$H_1: \beta_0 \neq 0 \quad \text{but not necessarily practically significant}$$

### confidence intervals for $\beta_0$

$$\text{RECALL: } \hat{\beta}_0 - \beta_0 \sim t_{n-(p+1)}$$

$$\text{Hence, } P(-t_{n-(p+1)} \leq \hat{\beta}_0 - \beta_0 \leq t_{n-(p+1)}) = 1 - \alpha.$$

$\hookrightarrow$  with prob.  $1 - \alpha$ , i.e. 100(1 -  $\alpha$ )% confidence:

$$\text{two-sided CI: } \hat{\beta}_0 - t_{n-(p+1)}^{\alpha/2} \text{SE}(\hat{\beta}_0) \leq \beta_0 \leq \hat{\beta}_0 + t_{n-(p+1)}^{\alpha/2} \text{SE}(\hat{\beta}_0) \quad \text{or} \quad (\hat{\beta}_0 - t_{n-(p+1)}^{\alpha/2} \text{SE}(\hat{\beta}_0), \hat{\beta}_0 + t_{n-(p+1)}^{\alpha/2} \text{SE}(\hat{\beta}_0)) \quad \text{enclosed}$$

result  $\hookrightarrow$  result

about 95% 2-sided CI

REMARKS:  $t_{\alpha/2}$  is the tail prob.:  $P(E = t_{\alpha/2}) = \alpha$ .

- choose  $\alpha = 0.05 \Rightarrow 100(1 - 0.05) = 95\%$ , i.e.

-  $t_{0.05}^{\alpha/2} \approx 2$ ; for large  $n-(p+1)$ ,  $t_{0.05}^{\alpha/2} \approx 1.96$

General Form: Estimate  $\pm$  (critical value) (SE(estimate))

$$\text{Corroborate } H_0: \beta_0 = 0 \quad \text{Given a } (1-\alpha) \text{ CI covers } \beta_0: 0 \text{ (Gd) p-value} > \alpha.$$

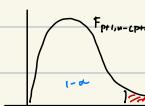
$$H_1: \beta_0 \neq 0$$

$$\text{so } \left| \frac{\hat{\beta}_0}{\text{SE}(\hat{\beta}_0)} \right| < t_{n-(p+1)}^{\alpha/2}$$

### Simultaneous confidence regions

Prop: With prob.  $1 - \alpha$ , i.e. confidence  $100(1 - \alpha)\%$

$$(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \leq (p+1)^2 F_{p+1, n-(p+1)}^{(1-\alpha)}$$



REMARKS: - Equivalently  $P[(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \leq (p+1)^2 F_{p+1, n-(p+1)}^{(1-\alpha)}] = 1 - \alpha$

$$\text{GIVEN: } \frac{1}{p+1} (\hat{\beta} - \beta)^T (\text{cov}(X^T X))^{-1} (\hat{\beta} - \beta) \sim F_{p+1, n-(p+1)} \quad \text{where } \text{var}(\hat{\beta}_{p+1}) = \sigma^2 (X^T X)^{-1}$$

$$\hookrightarrow \frac{1}{p+1} (\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \sim F_{p+1, n-(p+1)}$$

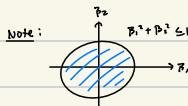
Prop: More generally, for any  $S \subseteq \{0, 1, \dots, p\}$

e.g.  $S = \{1, 2, 3\}$  or CI for  $(\beta_1, \beta_2, \beta_3)$ :

$$\frac{1}{|S|} (\hat{\beta} - \beta_S)^T [\text{var}(\hat{\beta}_S - \beta_S)]^{-1} (\hat{\beta}_S - \beta_S) \sim F_{|S|, n-|S|}$$

where  $|S|$ : cardinality of set  $S$

REMARKS: - when  $|S| = 1$  only one  $\beta_j$ , then this is equivalent to the 1-sided CI  $\beta_j$



NOTE: the correlation between predictors and coefficients are often different in sign

- two positively correlated predictors will attempt to perform the same job of explanation
- the more work one does, the less the other needs to do and hence a negative coefficient correlation

Proof (1): Given a single  $\beta_j$ :  $\frac{\hat{\beta}_j - \beta_j}{\text{SE}(\hat{\beta}_j)} \sim t_{n-(p+1)}$

$$\Rightarrow (\hat{\beta}_j - \beta_j)^2 = (\hat{\beta}_j - \beta_j)^T \frac{1}{\text{SE}^2(\hat{\beta}_j)} \sim F_{1, n-(p+1)}$$

$\hookrightarrow (\hat{\beta}_j - \beta_j)(\text{var}(\hat{\beta}_j))^{-1} (\hat{\beta}_j - \beta_j) \sim F_{1, n-(p+1)}$

Now consider  $\hat{\beta}_{p+1}$ : s.t.  $\hat{\beta}_{p+1} \sim N(0, \sigma^2 (X^T X)^{-1})$

Generalizing (1) to the multivariate case:

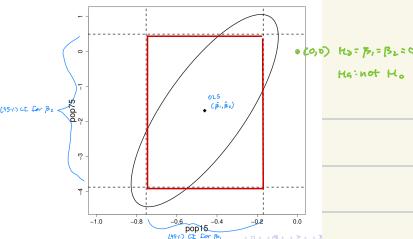
$$\frac{1}{p+1} (\hat{\beta} - \beta)^T (\text{cov}(X^T X))^{-1} (\hat{\beta} - \beta) \sim F_{p+1, n-(p+1)}$$

$\uparrow$

$\text{var}(\hat{\beta}_{p+1})$

$$\text{thus, } \frac{1}{p+1} (\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \sim F_{p+1, n-(p+1)}$$

use ellipsoid as result of stretching/shrinking by variance (and hence correlation) correction factor



Remark: simultaneous CI region more complex than intersection of respective 95% CI's (seen in red)  
→ Elliptical region accounts for simultaneous region; further details are beyond scope of class

- $(0, 0)$  outside of confidence region  $\Rightarrow$  p-value < 0.05
- direction (standardized) ellipse influenced by pairwise (column) correlation (vertical/horizontal shape)  $\Leftrightarrow$  orthogonal predictors

\* Better to consider joint CI's when possible - especially when  $\beta$ 's are heavily correlated

### CI's for Prediction

(Training) Data:  $(X_i, Y_i), i=1, \dots, n$ .

Given new predictors,  $X_0$ , what is the predicted response?

i.e. given  $Y_0 = X_0^T \beta + \epsilon_0$ , find  $\hat{Y}_0 = X_0^T \beta$ .

two types of predictions

1) Prediction ( $t_{n+1}$ ) of a future observation:  $= \hat{Y}_0 \pm t_{n+1}^{(1-\alpha)}$   $\hat{\sigma} \sqrt{X_0^T C X_0^{-1} X_0}$

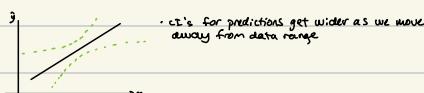
Remark: w.r.t.  $Y_0 = X_0^T \beta + \epsilon_0$

2) Prediction ( $t_{n+1}$ ) of the future mean response:  $= \hat{Y}_0 \pm t_{n+1}^{(1-\alpha)} \hat{\sigma} \sqrt{X_0^T C X_0^{-1} X_0}$

Remark:  $\text{var}(X_0^T \beta) = E(Y_0 | X_0)$

\* Prediction (future observation) intervals vs. confidence (future mean response) intervals

RESULT:  $C = \text{cov}(X_0, X_0) = X_0^T \beta$   
 $X_0 \in$  data frame  $(X_1, \dots, X_p)$   
Predict (residual,  $X_0$ , interval = "confidence")  
Predict (residual,  $X_0$ , interval = "prediction")



Prediction Band Plot:

### Extrapolation

- occurs when we try to predict the response for values of the predictor which lie outside the range of the original data; types:

- Quantitative  
→ check whether new  $X_0$  is within range of original data
- Qualitative  
→ check whether the new  $X_0$  is drawn from the same population from which the original sample was drawn

### Interpreting Parameter Estimates

$\beta_0$ : what does  $\beta_0, \beta_1, \dots, \beta_p$  mean?

$\beta_i$ : a unit increase in  $X_i$  wrt the other (column) predictors held constant will produce a change of  $\beta_i$  in the response  $y$ .

- may be problematic in practice  
e.g. idea of holding vars constant may not make sense for observational data (not under our control)

$$+ Y \sim X + X^2$$

Note: marginal effect of  $\beta_i$  on  $Y$  may change because of presence/absence of other predictors in model ( Simpson's Paradox - more on this later)

$\beta_0$ : expected value of response  $Y$  when all  $x$ 's = 0  
may not have literal sense depending on context

Proof (1). For a feature obs., where  $Y_0 = X_0^T \beta + \epsilon_0$ :

$$E(Y_0) = E(\hat{Y}_0) = X_0^T \beta \Leftrightarrow E(Y_0 - Y_0) = 0$$

$$\Rightarrow \text{var}(Y_0 - Y_0) = \text{var}(X_0^T \beta - X_0^T \beta - \epsilon_0)$$

$$= \text{var}(X_0^T (\beta - \beta)) + \text{var}(\epsilon_0) \text{ since } \beta \text{ and } \epsilon_0 \text{ are ind.}$$

(zero covariance)

$$= X_0^T \text{var}(\beta) X_0 + \sigma^2, \text{ higher dimension for } \text{var}(X_0) = n^2 \text{v}(X)$$

$$= \sigma^2 X_0^T C X_0 X_0 + \sigma^2$$

$$= \hat{\sigma}^2 (t_{n+1}^{(1-\alpha)})^2 X_0 + \sigma^2 \quad \text{and} \quad \frac{Y_0 - \bar{Y}}{\hat{\sigma} t_{n+1}^{(1-\alpha)}} \sim t_{n+1-p-1} \Rightarrow 1.$$

Proof (2). we know  $E(Y_0 | X_0) = X_0^T \beta$

Note that

$$E(\hat{Y}_0) = X_0^T \beta \quad \text{as above}$$

$$\text{var}(\hat{Y}_0 - X_0^T \beta) = \text{var}(X_0^T \beta - X_0^T \beta) = \hat{\sigma}^2 X_0^T C X_0 X_0$$

$$\text{and} \quad \frac{\hat{Y}_0 - X_0^T \beta}{\hat{\sigma} t_{n+1-p-1}} \sim t_{n+1-p-1} \Rightarrow 1.$$

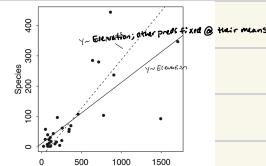


Figure 5.1 The fit for the simple model of just elevation as a predictor is shown as the solid line. The predicted response to observed values and the other four predictors held at their mean is shown as a dashed line.

**Effect Plot:** visualizing the mean of the model for a given predictor  
 + shows relationship between predictor and response for (logistic) model  
 + compute predicted responses for observed values of predictor, fixing the other predictors at their means

- observational Data
  - + causal conclusion problematic; steps to make stronger case for causality:
    - + try to include all relevant vars
    - + try a variety of models and see whether a similar effect is observed
    - + multiple studies under different conditions can help confirm a relationship
  - unmeasured lurking var Z may be the real cause and/or be a confounding variable

**Simpson's Paradox:** regression coefficient of a predictor may change sign when adding/removing another predictor;

$$\text{e.g. } E(Y|X, Z) = \beta_0 + \beta_1 X + \beta_2 Z$$

Data only have X and Y. F.i.  $Y \perp X$ :

$$E(Y|X) = E(\beta_0 + \beta_1 X + \beta_2 Z|X) = \beta_0 + \beta_1 X + \beta_2 E(Z|X)$$

Suppose  $\rightarrow X \perp Z$  (e.g. randomization)  $\Rightarrow E(Z|X) = E(Z)$ .

$$\Rightarrow E(Y|X) = [\beta_0 + \beta_2 E(Z)] + \beta_1 X$$

b) X and Z correlated (e.g. observational)

$$\text{and assume } Z = \gamma_0 + \gamma_1 X + \varepsilon \Rightarrow E(Z|X) = \gamma_0 + \gamma_1 X$$

$$\Rightarrow E(Y|X) = (\beta_0 + \beta_2 \gamma_1) + (\beta_1 + \beta_2 \gamma_1) X$$

### Experimental Design

Q: How can we establish causal conclusions?

- In **Designed Experiments**, our control over the experimental conditions allows us to make stronger conclusions from the analysis

\* Randomization can reduce the effect of unknown predictors not included in model

- + only reliable way to ensure data not unbalanced that favors either treatment/control
- + randomly assign experimental units to chosen values of the predictors
- + still doesn't necessarily ensure success; bad model, unrealistic controlled exp., etc.

\* Remark (Randomization): Let  $X_1$ : treatment (or not),  $X_2$ : unobserved (Data mean=0, var=1)

$$\Rightarrow \text{COV}(X_1, X_2) = 0 \Rightarrow \text{sample } X_1 \text{ and } X_2 \text{ are orthogonal}$$

$\rightarrow$  Motivates how randomization can achieve causal conclusions

### Orthogonality

- orthogonal predictors (by design) can also simplify the problem

- + allows us to more easily interpret the effect of one predictor w.r.t. another

Suppose we partition  $X_{n \times p+1} = [X_{1:n} \mid X_{2:n}]$  s.t.  $X_{1:n}^T X_{2:n} = 0$ .  
 $\underset{n \times n}{\underbrace{\qquad\qquad\qquad}} \quad \underset{n \times n}{\underbrace{\qquad\qquad\qquad}}$

$$\text{Now, } Y = X\beta + \varepsilon = X_{1:n}\beta_{1:n} + X_{2:n}\beta_{2:n} + \varepsilon$$

and  $(p+1) \times 2$

$$X^T X = \begin{pmatrix} X_{1:n}^T X_{1:n} & X_{1:n}^T X_{2:n} \\ X_{2:n}^T X_{1:n} & X_{2:n}^T X_{2:n} \end{pmatrix} = \begin{pmatrix} X_{1:n}^T X_{1:n} & 0 \\ 0 & X_{2:n}^T X_{2:n} \end{pmatrix}$$

which gives

$$\hat{\beta}_{1:n} = (X_{1:n}^T X_{1:n})^{-1} X_{1:n}^T Y \text{ and } \hat{\beta}_{2:n} = (X_{2:n}^T X_{2:n})^{-1} X_{2:n}^T Y$$

$$\hookrightarrow \hat{\beta}_1 = (X_{1:n}^T X_{1:n})^{-1} Y = \begin{pmatrix} (X_{1:n}^T X_{1:n})^{-1} & 0 \\ 0 & (X_{2:n}^T X_{2:n})^{-1} \end{pmatrix} \begin{pmatrix} X_{1:n}^T Y \\ X_{2:n}^T Y \end{pmatrix} Y$$

$$= \begin{pmatrix} (X_{1:n}^T X_{1:n})^{-1} X_{1:n}^T Y \\ (X_{2:n}^T X_{2:n})^{-1} X_{2:n}^T Y \end{pmatrix}.$$

\* note:  $\hat{\beta}_{1:n}$  will be the same regardless of whether  $X_{2:n}$  is in model or not

### Practical Difficulties

- Nonrandom Samples
- Choice of Range of Predictors
- Model Misspecification
- Publication and Experimenter Bias
- Practical and Statistical Significance

### Ch. 4: Diagnostics

Textbook	Lecture content
4.1: checking Error Assumptions	Checking Error Assumptions
4.2: Finding unusual observations	Practical Error Checking & Diagnostics
4.3: checking Structure of Model	Common Handling of Violated Assumptions

## Diagnostics

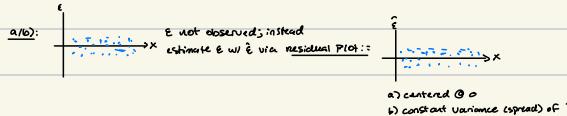
- checking error assumptions
- finding unusual points
- checking the model structure

### Checking error assumptions

Assume:

$$\begin{aligned} \text{E}(X) &\stackrel{\text{def}}{=} N(0, \sigma^2 I) \\ \text{a)} \quad \text{E}(\epsilon) &= 0 \\ \text{b)} \quad \text{Var}(\epsilon|X) &= \sigma^2 I \\ \text{c)} \quad \epsilon \text{ is } \text{indep. identically distributed, Normal} \end{aligned}$$

remark: typically check in this order



corollary: it can be shown that

$$\hat{\epsilon} = (I_n - H)\epsilon = (I_n - H)\hat{y}$$

$$\begin{aligned} \text{Proof: } \hat{\epsilon} &= Y - \hat{Y} = Y - HY, \quad H = X(X^T X)^{-1} X^T \Rightarrow \hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = HY \\ &\Rightarrow \hat{\epsilon} = (I_n - H)Y \end{aligned}$$

$$\text{and } HY = H(X\beta + \epsilon) = HX\beta + H\epsilon$$

$$= X\beta + H\epsilon, \text{ since } HX = X$$

$$\text{thus, } \hat{\epsilon} = \underbrace{X\beta}_{Y} - \underbrace{(X\beta + H\epsilon)}_{\hat{Y}} = \epsilon - H\epsilon$$

$$\Rightarrow \hat{\epsilon} = (I_n - H)\epsilon.$$

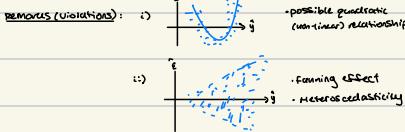
Residual Plot: Plot  $\hat{\epsilon}$  against  $x$  or any linear combination of  $x$ 's ( $\dots \hat{Y} = x^T \hat{\beta}$ ) to check for:

i) Non-linearity

→ expect a constant mean function equal to 0 if our model is correct, i.e.

ii) Homoscedasticity (Equal variance)

iii) Heteroscedasticity (Unequal variance)



note: we want variance of residuals to be constant (ii) across all levels of the independent variables

### Properties of Residuals

Treat  $\hat{\epsilon}$  as random variable

$$(1): E(\hat{\epsilon}|X) = 0 \Leftrightarrow E(Y|X) = X\hat{\beta}$$

$$(2): \text{Var}(\hat{\epsilon}|X) = \sigma^2(I_n - H) \text{ if } \text{Var}(\epsilon|X) = \sigma^2 I$$

$$(3): \text{cov}(\hat{\epsilon}|X) = 0$$

Treat  $\hat{\epsilon}$  as sample, then

$$i) \text{ mean of } \hat{\epsilon} = 0, \text{ i.e. } \bar{\hat{\epsilon}} = 0$$

$$\text{mean } \hat{\epsilon} = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i = 0$$

$$ii) \text{ Sample cov of } \hat{\epsilon} \text{ and } \hat{Y} = 0, \text{ i.e. } \text{cov}(\hat{\epsilon}, \hat{Y}) = 0$$

rem: (i) - (iii) hold if model assumptions satisfied;

(i), (ii) and (iii) always hold - even in spite of violated assumptions but not used for diagnosis

Q: Suppose SLR S.t.  $\hat{\epsilon} \sim \mathcal{N}$ .

$$\text{then } \hat{\beta}_1 = 0 \text{ and } \hat{\beta}_2 = \bar{y} - \bar{\hat{y}}$$

$$= \bar{y}$$

$$= \bar{x}$$

$$\text{Proof (1): } E(\hat{\epsilon}|X) = E((I_n - H)\epsilon|X) = (I_n - H)E(\epsilon|X) = 0, \quad \text{Proof (2): } \text{cov}(\hat{\epsilon}, \hat{\epsilon}) = \frac{1}{n} \sum_{i=1}^n (\hat{\epsilon}_i - \bar{\hat{\epsilon}})(\hat{\epsilon}_i - \bar{\hat{\epsilon}}) \quad \text{where } \bar{\hat{\epsilon}} = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i = 0 \text{ by (i)}$$

$$\text{Proof (3): } \text{var}(\hat{\epsilon}|X) = \text{Var}((I_n - H)\epsilon|X)$$

$$= (I_n - H)\text{Var}(\epsilon|X)(I_n - H)^T, \text{ since } (I_n - H) \text{ is symmetric}$$

$$= \sigma^2(I_n - H)(I_n - H)^T$$

$$= \sigma^2(I_n - 2H + H^2)$$

$$= \sigma^2(I_n - H), \text{ since } H^2 = H.$$

$$\text{and } \bar{Y} = \frac{1}{n} \sum_{i=1}^n \hat{Y}_i$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \bar{Y}_i - (\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i) \bar{Y}$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \bar{Y}_i - \bar{\hat{\epsilon}} \bar{Y} \quad \text{by (i)}$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \bar{Y}$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \bar{Y} \quad \text{by (i)}$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \bar{Y} \quad \text{by (i)}$$

$$= 0 \quad \text{by normal eq. } X^T \hat{\beta} = 0 \Leftrightarrow X^T \bar{\hat{\epsilon}} = 0$$

$$= 0 \in \mathbb{R}^{n \times 1}$$

$$\text{Proof (4): } \text{cov}(\hat{\epsilon}, \hat{Y}|X) = \text{cov}((I_n - H)\epsilon, HY|X)$$

$$= (I_n - H)\text{Var}(Y|X)H^T$$

$$= \sigma^2(I_n - H)H^T$$

$$= \sigma^2(I_n - H^2)$$

$$= 0.$$

## Assumption checking (Practically)

1) zero mean function assumption (linearity),

e.g. fit  $\hat{y} \sim \hat{y} + \hat{y}^2$  and check if  $\hat{y}^2$  is significant

2) Fit  $\hat{y} \sim \hat{y}$  and check (test) whether slope term is significant (constant variance)

(evidence for non-constant variance (heteroscedasticity))

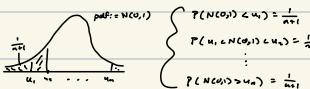
3) Q-Q Plot (Normality)

a) Sort the residuals  $\hat{e}_{(1)} \leq \hat{e}_{(2)} \leq \dots \leq \hat{e}_{(n)}$

b) Compute  $u_i = \Phi^{-1}\left(\frac{i}{n+1}\right), i=1, \dots, n$  where  $\Phi$ : cdf of  $N(0,1)$  S.t.

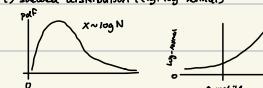
↳ theoretical normal quantiles

$$\Phi(x) = P(N(0,1) \leq x)$$



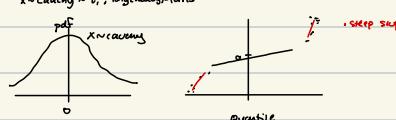
c) Plot  $\hat{e}_{(i)}$  against  $u_i$ :

Ex (non-normality): i) skewed distribution (e.g. log-normal)



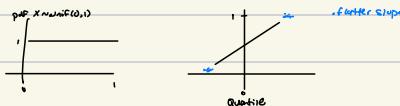
ii) long-tailed dist. (e.g. cauchy)

$X \sim \text{cauchy} \sim t_1$ ; long/heavy-tails



iii) short-tailed dist. (e.g. uniform)

$X \sim \text{unif}$ ; short/light-tail



Shapiro-Wilk test for normality:  $H_0$ : residuals are normal

• not very helpful or used in practice (Q-Q plot better)

• small  $n$ : little power

• Large  $n$ : non-normality is less important; non-normality mitigated by large  $n$

## 4) Correlated Errors

In temporally related data:

• Plot  $\hat{e}_t$  against time

• Plot  $\hat{e}_t$  against  $\hat{e}_{t-1}$

• time series analysis may be more appropriate

Remark: If there is no temporal relationship or other ordering in the variables, checking independence is more challenging

Durbin-Watson test for correlated errors:  $H_0$ : uncorrelated errors

- can create an Index Plot: + look for successive patterns as evidence for correlated residuals

## 5) Finding Unusual Points

outliers: do not fit the model well

influential points: affect the fit of the model substantially

Note: a point can be none, one, or both of these

Leverage: Recall  $H = X(X^T)^{-1}X^T$

$$\boxed{\text{Leverage of point } i = h_{ii} = h_{ii}^2}$$

$h_{ii}$  only depends on  $X$

$$\rightarrow \text{var}(e_i) = \sigma^2(1-h_{ii}) \quad ; \text{ recall: var}(e_i) = \sigma^2(I - H)$$

$\sum h_{ii} \leq n-1$ ; By definition,  $h_{ii} = x_i^T(X^TX)^{-1}x_i$  where  $x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$

$$x_i^T = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$\sigma^2(I - H)$$

$$h_{ii} \leq 1$$

Let  $x_i^H = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$ . Then,

$$h_{ii} = \frac{1}{n} + (x_i^H - \bar{x}^H)^T(X^TX)^{-1}(x_i^H - \bar{x}^H)$$

$$\text{where } \bar{x}^H = \frac{1}{n} \sum x_i^H, X^H = \begin{pmatrix} (x_1^H)^T \\ (x_2^H)^T \\ \vdots \\ (x_n^H)^T \end{pmatrix}$$

Suppose  $x_i^H = R_{ip}$ , i.e.

$$\rightarrow \sum h_{ii} = p+1; \quad \sum h_{ii}h_{jj} = \text{tr}(H) = \text{tr}(X(X^T)^{-1}X^T)$$

$$= \text{tr}((X^T)^{-1}X^T)$$

$$= \text{tr}(I_{pp}) = p+1.$$

- Rule of thumb: leverages greater than  $\frac{2(p+1)}{n}$  are considered 'high'

Note:  $t_{ij} = \frac{1}{n} E h_{ij} = \frac{1}{n} (p+1)$

Remark: If  $h_{ii} > 1$ , then

$$\text{Var}(e_i) = \sigma^2(1-h_{ii}) \rightarrow 0 \Rightarrow e_i \approx 0$$

#### Half-Normal Plot

- Sort  $h_{11}, h_{22}, \dots, h_{nn}$
- Compute  $u_i = \Phi^{-1}\left(\frac{h_{(i)}}{n+1}\right)$
- Plot  $h_{ii}$  against  $u_i$
- $h_{ii} > \frac{2(p+1)}{n}$ ; points above considered to have 'high' leverage

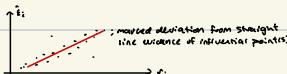


#### Studentized Residuals

Since  $\text{Var}(e_i) = \sigma^2(1-h_{ii})$ , define

$$\text{(internally) studentized residuals: } t_i = \frac{e_i}{\hat{\sigma} \sqrt{1-h_{ii}}} \quad \text{where } \hat{\sigma} \sqrt{1-h_{ii}} = \hat{\sigma} \sqrt{V_i}$$

- Better to use studentized residuals for diagnostic plots (e.g. Q-Q plot and testing constant variance)
- Practically, results are often similar (as compared w/ using  $e_i$ )



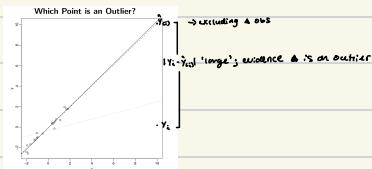
#### Outliers

Q: How do we deal with both truly unusual points and large residuals?

• Exclude point  $i$ , recompute  $\hat{y}_{(i)}$  and hence  $\hat{g}_{(i)} = X_i^T \hat{\beta}_{(i)}$

• If  $|y_i - \hat{g}_{(i)}|$  is large, then obs.  $i$  is an outlier

Q: How large  $i$  is large



#### (Externally) Studentized Residuals

It turns out

$$t_i = \frac{y_i - \hat{g}_{(i)}}{\hat{\sigma} \sqrt{1-h_{ii}}} \quad (1) \quad \text{where } \hat{\sigma} \sqrt{1-h_{ii}} \text{ same as prediction interval for a new obs. (Eqn. 4)}$$

$$= r_i \left( \frac{n-p+1-i}{n-p+1-i-1} \right)^{1/2} \quad (2)$$

$$\sim t_{n-(p+1)-1} \text{ under } H_0: \text{obs. } i \text{ not outlier} \quad (3)$$

#### Multiple Hypothesis Tests

• If  $|t_{ii}|$  is too large (small p-values), reject  $H_0$  and conclude obs.  $i$  is an outlier

• For each obs.  $i$ , compare  $|t_{ii}|$  with  $t_{n-(p+1)-1}^{1/2}$

Note: when (simultaneously) conducting more than one hypotheses,

# we will reject too many points

#### Bonferroni Correction

Type I Error (rate) =  $P_{H_0}(\text{reject at least one test})$

$$\leq \sum_{i=1}^n P_{H_0}(\text{reject test } i)$$

$$= n\alpha$$

thus, introduce the Bonferroni correction: test each hypothesis at level  $\alpha/n$

#### Remarks on outliers

- Two/more outliers could hide each other
- Cluster of outliers: consider using robust methods
- Examine the context - what could it mean?
  - + data entry errors
  - + lurking variables
  - + something gone wrong, e.g. Gradient credit card use
  - + a new unknown effect
  - + some patterns just have exceptions

### Influential Points

DEF: one whose removal from the dataset would cause a large change in the fit



- Measure the influence via,

- +  $\hat{\beta}_j = \hat{\beta}_{(j)}$
- + old using old data with obs.
- +  $K(\hat{\beta} - \hat{\beta}_{(j)}) = \hat{y} - \hat{y}_{(j)}$
- +  $|e_{(j)}|$

### Cook's Distance

$$\text{Cook statistic: } D_i = \frac{(y - \hat{y}_{(i)})^2 (\hat{y} - \hat{y}_{(i)})}{\text{const. } \hat{\sigma}^2}$$

$$= \frac{1}{p+1} \cdot \frac{t_i^2}{1-t_i^2} \quad \left\{ \begin{array}{l} t_i \geq 3 \text{ standardized (studentized) residuals} \\ t_i = \text{externally studentized residuals} \end{array} \right.$$

+ combination of residual and leverage effect

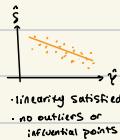
+  $D_i \approx 1$  is a reference value for potential influential points

### Checking the Structure of the Model

- Plot  $\hat{\epsilon}$  against  $\hat{y}$  and  $x_j$  - but other predictors impact the relationship;  
want to isolate the effect of  $x_j$  on  $y$  (conditional on the other predictors)

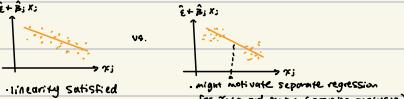
consider:

- i) Partial Regression Plots (Added Variable Plots)
- ii) Regress  $y$  on all  $x$  except  $x_j$ ; get residuals  $\hat{\epsilon}$
- iii) Plot  $\hat{\epsilon}$  against  $\hat{x}_j$ 
  - slope is  $\hat{\beta}_j$ ; look for non-linearity, outliers, and influential points
  - note: same as  $\hat{\beta}_j$ ; when fitting  $y$  on all  $x$ 's;
  - LSL-UCL help add res more diagnostics



### (z): Partial Residual Plots : Plot $\hat{\epsilon} + \hat{\beta}_j x_j$ against $x_j$

- B: where does this come from?  
 $A: Y = \sum_{j \neq j} x_j \hat{\beta}_j + x_j \hat{\beta}_j + \hat{\epsilon}$   
 $\Rightarrow Y = \sum_{j \neq j} x_j \hat{\beta}_j + x_j \hat{\beta}_j + \hat{\epsilon}$



### Summary of Diagnostics

- Just fitting a model is not good enough
- Graphical diagnostics are more informative but also more subjective
- Diagnostics often suggest a change in the model and then the whole process is repeated
- time-consuming but necessary

### What to do w/ violated assumptions

- Heteroscedasticity
  - + nothing
  - + use other methods for inference, e.g. bootstrap
- Weighted Least Squares (Ch. 9)
  - + transformation of the response (Ch. 9)
- Nonlinearity: change the model
  - { add new predictors, e.g.  $x_1^2, x_1^3, \dots$
  - + transformation (Ch. 9)
- Non-Normality
  - + nothing
  - + especially for large n use CLT and apply z-table for short-tailed dist - not heavy tail dist (e.g. Cauchy)
  - + Transformation of the response (skewed errors)
  - + Robust methods (long-tailed dist.)
  - + Inference based on other distributions
- Correlated Errors
  - + Generalized Least Squares (Ch. 6)

## Ch. 5: Problems with the Predictors

Textbook	Lecture content
S.1: Errors in the Predictors	(MSMSE) Error in Predictors
S.2: changes of scale	change of scale
S.3: collinearity	Standardizing variables Collinearity

- Errors in Predictors: suppose  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ ,  
 then if  $x_i \in X_i + \delta_i$ ,  $\delta_i$  (measurement) error  
 . measurement error in X (preds) usually would lead to bias  
 . if  $\text{var}(\delta) \gg \text{var}(x)$ , then the bias is small and can practically be ignored

- Change of Scale: • predictors of similar magnitude are easier to compare

- numerical stability
- can aid interpretation

- Consequences: • rescaling  $x_i$  leaves the t and F tests and  $\hat{\sigma}^2$  and  $R^2$  unchanged.

Suppose  $x_j \rightarrow \frac{x_j + a}{b}$ , j-th predictor

while  $\hat{\beta}_j \rightarrow b\hat{\beta}_j$

$$\sqrt{\text{var}(\hat{\beta}_j)} = \hat{\sigma}(\hat{\beta}_j) \rightarrow b\hat{\sigma}(\hat{\beta}_j) = \sqrt{b^2 \text{var}(\hat{\beta}_j)}$$

$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$

$$\text{with } x_j \rightarrow \frac{x_j + a}{b} \text{ and } \beta_j x_j = (\beta_j/b) \frac{x_j + a}{b} - \beta_j a$$

$\Rightarrow \beta_j \rightarrow \beta_j/b$  and  $\beta_0 \rightarrow \beta_0 - \beta_j a$ , other  $\beta_j$ 's,  $\epsilon$ 's unchanged

- rescaling y leaves the t and F tests and  $R^2$  unchanged but both  $\hat{\sigma}$  and  $\hat{\beta}$  rescaled by  $b$ ;  $\hat{\beta}$  is both shifted by  $a$  and rescaled by  $b$ :

Suppose  $y \rightarrow \frac{y+a}{b}$ .

$$\text{then } \frac{y+a}{b} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon + a$$

$$\Rightarrow \beta_0 \rightarrow \frac{\beta_0 + a}{b}, \beta_1 \rightarrow \frac{\beta_1}{b}, \dots, \beta_p \rightarrow \frac{\beta_p}{b}, \epsilon \rightarrow \frac{\epsilon}{b}$$

- Standardizing Variables: • convert all variables to standard units (mean=0, sd=1), i.e.

$$\text{for } i=1, \dots, n, j=1, \dots, p, x_{ij} \rightarrow \frac{x_{ij} - \bar{x}_j}{\text{sd of } x_{ij}}$$

$$\text{for } i=1, \dots, n, y_i \rightarrow \frac{y_i - \bar{y}}{\text{sd of } y}$$

- can compare coefficients directly

- helps numerical stability, e.g.

$(X^T X)^{-1}$  computation

- interpretation is harder

note: Suppose  $X$  and  $y$  standardized. Then,

$$\tilde{y} = \hat{\beta}_0 + \hat{\beta}_1 \tilde{x}_1 + \dots + \hat{\beta}_p \tilde{x}_p$$

but all  $\tilde{x}_1, \dots, \tilde{x}_p$  are zero

$$\Rightarrow \hat{\beta}_0 = 0.$$

- Collinearity:  $X^T X$  close to singular (not invertible)

- Cause: some columns are (almost) linear combos of others

### Detection

- correlation matrix: large pairwise correlation  
 → may not work for multicollinearity involving more than two predictors, e.g.  $x_1 \approx x_2 + x_3$
- regress  $x_i$  on other preds - get  $R^2_i$ :  
 →  $R^2$  close to 1 (e.g.  $R^2 \geq 0.9$ ) indicates a problem
- cond. no. number of  $X^T X$ :  $= K = \frac{\lambda_1}{\lambda_{p+1}}$

where  $\lambda_1$ : largest eigenvalue of  $X^T X$

$\lambda_{p+1}$ : smallest eigenval of  $X^T X$

e.g. suppose  $X^T X = \begin{pmatrix} 1 & \dots \\ \vdots & \ddots & 1 \end{pmatrix} \Rightarrow \lambda_1 = \sqrt{\lambda_{p+1}} = 1$

$$X^T X = \begin{pmatrix} 1 & \dots & 0 \\ \vdots & \ddots & 0 \\ 0 & \dots & 0 \end{pmatrix} \Rightarrow \lambda_1 = \sqrt{\lambda_{p+1}} = \sqrt{100} = 10$$

\* large  $K$  indicates (evidence) of collinearity

→ practically,  $K \geq 20$  strong (presence) of evidence

Consequences: - imprecise estimate of  $\beta$ , i.e.

$\text{se}(\hat{\beta}_j)$  large

- t-test fails to reveal significant pred:

$$\text{se}(\hat{\beta}_j) \text{ large} \Rightarrow t = \frac{\hat{\beta}_j}{\text{se}(\hat{\beta}_j)} \rightarrow 0 \rightarrow \text{large p-value} \rightarrow \text{FTR Ho}$$

- sensitivity to measurement errors

- numerical instability

Variance Inflation Factor: Let  $S_{xy} = \sum_i (x_{ij} - \bar{x}_j)^2$ , then, it can be shown that

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \left( \frac{1}{1 - R_j^2} \right) \frac{1}{S_{yy}}$$

- if  $x_j$  is independent (orthogonal) from other predictors,  $\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{S_{yy}}$

- if  $R_j^2 \rightarrow 1$ ,  $\text{Var}(\hat{\beta}_j) \rightarrow \infty$

$$\text{and we define: } \text{VIF} := \frac{1}{1 - R_j^2}$$

+ (large) VIF indicates a problem  
e.g.  $\text{VIF} > 10 \text{ or } R_j^2 > 0.9$

How to resolve collinearity: - if you mostly care about inference, drop highly correlated predictors

- variable selection may be used

- if interpretation is important and you must keep all predictors, do not use least squares; use some other estimation method, e.g. ridge regression

## ch. 6: Problems with the Errors

Textbook	Lecture content
6.1: Generalized Least Squares	weighted least squares
6.2: Weighted Least Squares	Generalized Least Squares
6.3: Testing for lack of Fit	Robust Regression
6.4: Robust Regression	

### What can go wrong w/ the errors

Assume  $E(\epsilon | x) = 0$ ,  $\text{Var}(\epsilon | x) = \sigma^2 I$ , i.e.

↳ unequal variances:  $\text{Var}(\epsilon_i) \neq \text{Var}(\epsilon_j)$

↳ correlated:  $\text{corr}(\epsilon_i, \epsilon_j) \neq \text{cor}(\epsilon_i, \epsilon_j)$

↳ heavy-tailed

### Weighted Least Squares

Errors uncorrelated, but unequal variance, i.e.

$$\epsilon \sim N(0, \sigma^2 W^{-1}) \quad \text{where } \sigma^2 W^{-1} = \begin{pmatrix} \sigma^2 w_{11} & & \\ & \ddots & \\ & & \sigma^2 w_{nn} \end{pmatrix}_{n \times n}$$

and  $W^{-1} = \text{diag}\left(\frac{1}{w_1}, \dots, \frac{1}{w_n}\right)$ ,  $w = \left(\frac{w_1}{\sigma^2}, \dots, \frac{w_n}{\sigma^2}\right)_{n \times n}$

Remark: If  $w_i \sim \text{Gamma}(1, \sigma^2) \Rightarrow \text{Var}(w_i) = \sigma^2 I$  (OLS)

↳ Error variance proportional to a predictor, i.e.

$$\text{Var}(\epsilon_i) \propto x_{i1} \Rightarrow w_i = x_{i1}^{-1}$$

↳ suppose  $y_i$  is the average of  $n_i$  observations. Then  $w_i \propto n_i$

$$\text{Proof: } \text{Var}(y_i) = \text{Var}\left(\frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \epsilon_j\right)$$

assuming  $x_{ij}$  homoskedastic,  $\text{Var}(x_{ij} \epsilon_j) = \sigma^2$

$$= \frac{\sigma^2}{n_i} \Rightarrow w_i = n_i^{-1}$$

### Estimates

$$\text{Under (WLS) model: } y_i = \beta^T x_i + \epsilon_i, i=1, \dots, n, \text{Var}(\epsilon_i) = \sigma^2 w_i$$

( $\Leftrightarrow y = X \beta + \epsilon$ ,  $\epsilon \sim N(0, \sigma^2 W^{-1})$ )

Note that  $\text{Var}(W^{1/2} \epsilon_i) = \text{Var}(W_i^{1/2} \epsilon_i) = \sqrt{w_i} \cdot \text{Var}(\epsilon_i) = w_i \frac{\sigma^2}{w_i} = \sigma^2$

Now, assume the following

$$\text{transformation: } \begin{cases} x_i \rightarrow \sqrt{w_i} x_i \\ z_i \rightarrow \sqrt{w_i} x_i \end{cases}$$

$$\text{thus, } \text{Var}(y_i) = \beta^T \sqrt{w_i} x_i + \text{Var}(\epsilon_i)$$

$$z_i = \beta^T z_i + \epsilon_i^*$$

which satisfy OLS assumptions.

$$\Leftrightarrow y_i^* = x_i^T \beta^* + \epsilon_i^*$$

$$\text{with } y_i^* \text{ with } \text{Var}(y_i^*) = \sigma^2 z_i$$

$$\text{Thus, } \hat{\beta} = ((X^T W X)^{-1}) (X^T W Y) \\ = (X^T W^T W^T X)^{-1} X^T W^T W Y \\ = (X^T W X)^{-1} X^T W Y \quad (1).$$

Similarly,  $\text{var}(\hat{\beta}) = \sigma^2 ((X^T W X)^{-1})$   
 $= \sigma^2 (X^T W X)^{-1} \quad (2)$

Letting  $\hat{E} = Y - X\hat{\beta}$ , so,

$$\hat{E}^T Y = Y^T \hat{E} = W^T (Y - X\hat{\beta}) = W^T \hat{W} \hat{E}.$$

$$\text{RSS} = \hat{E}^T \hat{E} = \hat{E}^T W \hat{E}.$$

Thus,  $\hat{\sigma}^2 = \frac{1}{n-p+1} \hat{E}^T W \hat{E} \quad (3)$ ; thus:

$\begin{aligned} \hat{\beta} &= (X^T W X)^{-1} X^T W Y \\ \text{estimates: } &\left\{ \begin{array}{l} \hat{\beta} = (X^T W X)^{-1} X^T W Y \\ \text{var}(\hat{\beta}) = (X^T W X)^{-1} \sigma^2 \\ \hat{\sigma}^2 = \frac{\hat{E}^T W \hat{E}}{n-p+1} \end{array} \right. \\ \hat{\sigma} &= \sqrt{\frac{\hat{E}^T W \hat{E}}{n-p+1}} \end{aligned} \right.$
---

REMARKS: non-null weights can be used to indicate that different obs. have different variances  
 + otherwise, unbalanced groups can cause biased estimates

### Generalized Least Squares

Errors (potentially) correlated and unequal variance, i.e.

In general,

$E \sim N(0, \sigma^2 \Sigma)$ where $\Sigma$ is a general covariance matrix that satisfies $\left\{ \begin{array}{l} \text{Positive Def} \\ \text{Symmetric} \end{array} \right.$
--

e.g.  $\Sigma^{-1/2} E \sim N(0, \sigma^2 I)$  since  $\text{var}(\Sigma^{-1/2} E) = \Sigma^{-1/2} \sigma^2 \Sigma \Sigma^{-1/2} = \sigma^2 I_n$

Practically,  $\Sigma^{-1/2}$  is not uniquely defined. Instead,

<p>choose <math>\Sigma = S S^T</math></p> <p>where <math>S</math>: lower triangular matrix (the Cholesky decomposition)    and</p> $S^T E \sim N(0, \sigma^2 I) \Rightarrow \text{var}(S^T E) = \sigma^2 I_n$
---

So, we have the following

<p>transformation: <math>\begin{cases} y \rightarrow S^T y \\ x \rightarrow S^T x \\ \epsilon \rightarrow S^T \epsilon \end{cases} \Rightarrow Y = X\hat{\beta} + \epsilon, \text{ var}(\epsilon) = \sigma^2 \Sigma</math></p> <p><math>\downarrow</math></p> <p><math>S^T Y = S^T X\hat{\beta} + S^T \epsilon</math></p> <p><math>\underbrace{y^T}_{y^T S^T} \underbrace{x^T}_{x^T S^T} \underbrace{\epsilon^T}_{\epsilon^T S^T}, \text{ var}(\epsilon) = \sigma^2 \Sigma, \text{i.e. OLS}</math></p>
--

With

<p>estimates: <math>\left\{ \begin{array}{l} \hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y \\ \text{var}(\hat{\beta}) = (X^T \Sigma^{-1} X)^{-1} \sigma^2 \\ \hat{\sigma}^2 = \frac{\hat{E}^T \Sigma \hat{E}}{n-p+1} \end{array} \right.</math></p>
--

REMARKS: • WLS is a special case of GLS with  $\Sigma = W^T W$   
 • Practically, we estimate  $\Sigma$  from data

considering correlated Errors: both collected over time: errors could be correlated

• one of the simplest correction structures over time is the

autoregressive model:  $E_{t+1} = P E_t + \varepsilon_t$

<p>assuming <math>\left\{ \begin{array}{l} \text{(AR(1)) i.e. today's statistical error only depends on yesterday's} \\ \text{stationary, i.e. } \varepsilon_t \text{ has same mean and variance, } t=1, \dots, n \\ \varepsilon_t \stackrel{iid}{\sim} N(0, \sigma^2) \end{array} \right.</math></p>
---

<p>then: <math>E(\varepsilon_{t+1}) = P E_t + E(\varepsilon_t) = 0 \quad \text{var}(\varepsilon_{t+1}) = P^2 \text{var}(\varepsilon_t) + \text{var}(\varepsilon_t)</math></p> <p>and <math>E(\varepsilon_{t+1}) = 0 \quad \text{var}(\varepsilon_{t+1}) = \sigma^2</math></p> <p>and <math>E(\varepsilon_{t+1}) = E(\varepsilon_t) = 0 \quad \sigma^2 = \sigma^2 = \gamma^2</math></p> <p>and <math>P = P(\alpha) \Rightarrow \text{det}(P) = \alpha</math></p>
---

<p><math>\Rightarrow E(\varepsilon_{t+1}) = 0 \quad \Rightarrow \text{var}(\varepsilon_{t+1}) = \frac{\sigma^2}{1-\alpha^2}</math></p> <p>and <math>\text{cov}(\varepsilon_t, \varepsilon_{t+1}) = \text{cov}(\varepsilon_t, P\varepsilon_t) = P^{1-t+1} \sigma^2</math></p>
--

so

$$\Sigma = \begin{pmatrix} \frac{\sigma^2}{1-\alpha^2} & & \\ & \ddots & \\ & & \frac{\sigma^2}{1-\alpha^2} \end{pmatrix}$$

Notice:  $\Sigma$  depends on unknown  $\alpha$ ,  $P$  parameters

Prob. can test  $H_0: P = 0$   
 H: Prob. evidence for correlated errors, i.e.  
 $\text{corr}(\varepsilon_t, \varepsilon_{t+1}) \neq 0$   
 DD for  $P$  does not contain  $\alpha$

### Robust Regression

Main concern: heavy-tailed error distribution

1) M-estimation (A): what we cover

2) Least-trimmed estimation (in textbook if interested)

<p>M-Estimation: Find <math>\hat{\beta}</math> to minimize <math>\sum_{i=1}^n L(Y_i - X_i \hat{\beta})</math></p>
---

REMARKS: •  $L(\cdot)$  is the loss function, possible loss functions:

$$+ L(x) = x^2: \text{least squares (LS)}$$

Robust Estimation: less sensitive to large (residual) values

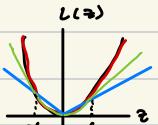
should plot distribution of residuals

+  $L(\beta) = \sum_i |y_i - \hat{y}_i|$  least absolute regression (LAD);  
equivalently called median regression

$$\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\Rightarrow \text{median}(Y|X) = X^T \beta$$

note: if dist. of  $\epsilon_i$  is symmetric, then  
median = mean



Remark: in OLS:  $\min_{\beta} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \Rightarrow \hat{\beta} = \bar{y}$

For LAD:  $\min_{\beta} \sum_{i=1}^n |y_i - \hat{y}_i| \Rightarrow \hat{\beta} = \text{median}(Y|X)$

+ Huber's method: =  $\begin{cases} \hat{\beta} = c, & \text{if } |c| \leq 1 \\ \text{OLS}, & \text{otherwise} \end{cases}$

$c$  should be a robust estimate of  $\sigma$ , e.g. median of  $|z_i|$

## Ch. 7: Transformation

### Textbook

7.1: transforming the response  
7.2: transforming the predictors

### Lecture content

Transforming the response (Box-Cox method)  
Transforming the Predictors

#### Outline

##### Transforming the response

- Box-Cox method

##### Transforming the predictors

- Polynomials

- Regression splines

#### Reasons to try transformations

- Nonlinearity

Suppose  $y=X$ : nonlinear but

$g(z)=X$ : linear for  $g$ : transformation function

##### Model misspecification

- May improve fit

- Incorporate a physical law (or some other known relationship?)

e.g.  $w \rightarrow \text{Area}$     $\text{Area} = L \cdot W$   
 $\log \text{Area} = \log L + \log W$

#### Box-Cox Method

Transformation of the response  $y \rightarrow g_\lambda(y)$

A family of transformations indexed by  $\lambda$  when  $y > 0$ :

$$g_\lambda(y) = \begin{cases} \frac{y^{\lambda-1}}{\lambda-1}, & \lambda \neq 0 \\ \ln(y), & \lambda = 0 \end{cases}$$

Remarks:

$\lambda=1 \rightarrow$  no transformation  
 $\lim_{\lambda \rightarrow 0} \frac{y^{\lambda-1}}{\lambda-1} = \log(y) = \ln(y)$

Prop. For  $\lambda \neq 0$ ,  $\frac{y^{\lambda-1}}{\lambda} \rightarrow [G_m(y)]^{1/\lambda}$  makes the unit of  $g_\lambda(y)$  the same for different  $\lambda$

$$\text{where } G_m(y) = \exp \sum_{i=1}^n \log(y_i) \bar{y} := \text{geom mean of } y_1, \dots, y_n$$

Goal: Find  $\lambda$  s.t.  $g_\lambda(y) \sim X$  follows linear model, i.e.

we minimize  $\text{RSS}_\lambda$  from regression  $g_\lambda(y) \sim X$ .

Equivalently, we can instead

compute likelihood of the data using the normal assumption

for any given  $\lambda$ .

then,

choose  $\lambda$  to maximize  $L(\lambda) = -\frac{n}{2} \ln \left( \frac{\text{RSS}_\lambda}{n} \right)$  (quadratic negative function)

compute confidence intervals for  $\lambda$  using asymptotic dist. of the likelihood

Remarks:

- may not choose  $\lambda$  that exactly maximizes  $L(\lambda)$  (i.e.  $\min_n \text{RSS}_\lambda$ ), but instead choose one that is easily interpreted, e.g.

$$\lambda = -2, -1, 0, \pm \frac{1}{2}, \pm \frac{1}{3}, \dots$$

sensitivity to outliers, e.g.

$\lambda = 5$  is extreme; practically,  $\lambda \in [-2, 2]$  is often used  
mainly due to outliers

- If some  $y_i \leq 0$ , can add constant  
i.e. transform  $y$  s.t.  $Y \geq 0$  e.g.  $\min(Y) = -1 \Rightarrow Y^* = Y + 1$  is now valid for BC-transformation
  - Transformations of proportions, counts - elect to use generalized linear models (see later)
  - If  $y_i$ 's are proportions (i.e.  $0 < y_i < 1$ ), consider logit transformation:  $\ln\left(\frac{y}{1-y}\right)$
- Note: log transformation is often used in practice if the range of a variable ( $y$  or  $x$ ) is over several magnitudes (e.g. home)

### Transforming the Predictors

Given  $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$  for predictors  $\begin{cases} x_1 \\ \vdots \\ x_p \end{cases}$

Now suppose

$Y = \beta_0 + \beta_1 f_1(x) + \dots + \beta_p f_p(x) + \epsilon$  for (new) predictors  $\begin{cases} f_1(x) \\ \vdots \\ f_p(x) \end{cases}$

Basis functions:  $f_j(x)$ , e.g.  $\begin{cases} \text{Polynomials} \\ \text{Regression splines} \end{cases}$

#### Polynomials: one Predictor case

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon \quad (R: \ln(Y \times X + I(X^n) + \dots + I(X^m)))$$

① How to choose  $d$ :

Forward: keep adding terms until new term is not statistically significant

Backwards: start w/ large  $d$  & keep eliminating the non-significant highest order term

REMARKS (ISSUES):

- (i) as  $d \rightarrow \infty$ ,  $x_i^{d+1}$  and  $x_i^d$  may be highly correlated
- (ii) as  $d \rightarrow \infty$ ,  $x_i^d$  too large (if  $|x_i| > 1$ ) or too small (if  $|x_i| < 1$ )  $\Rightarrow (x_i^d x_i)^{-1}$  problematic
- usually don't eliminate lower order terms, even if not significant, when a higher order term is significant (hierarchy)
- transforming predictor (e.g.  $X_1 + X_2 \rightarrow 0$ ) can affect significance of higher-order terms w.r.t. predictor (X)
- recall: previously familiar w/ unpaired t-test'd for transformed preds (w.r.t. decisiveness linear terms, i.e. degree)

### Orthogonal Polynomials

Suppose  $p_1, p_2, \dots, p_n$  for original predictors.

Assume polynomial original predictors ( $d$  degrees),  $i=1, \dots, n$  s.t.  $x_i = \begin{pmatrix} x_{i,0}=1 \\ x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,d}=(d+1)x_i \end{pmatrix}$

Then, for numerical stability,  $V \in \mathbb{R}^{n \times d+1}$ :

new predictors:	$\begin{cases} z_{i,0} = 1 \text{ (intercept)} \\ z_{i,1} = a_1 + b_1 x_{i,1} \\ z_{i,2} = a_2 + b_2 x_{i,2} + c_2 x_{i,1}^2 \Rightarrow z_{i,2} = \begin{pmatrix} 1 \\ a_2 \\ b_2 \\ \vdots \\ c_2 \\ (d+1)x_i \end{pmatrix} \\ \vdots \\ z_{i,d} = a_d + b_d x_{i,d} + \dots + c_d x_{i,1}^{d-1} \end{cases}$
-----------------	---

$(R: \ln(Y \sim Poly(K, d), "raw" = F))$

note: "raw" = T returns higher-order terms for preds w/o orthogonal transform

where  $a_j, b_j, c_j \dots$  are chosen s.t.  $z_i^T z_j = 0$  (orthogonal) when  $j \neq i$ ; i.e.

$$\text{where } \tilde{z}_j = \begin{pmatrix} z_{j,0} \\ z_{j,1} \\ \vdots \\ z_{j,d} \end{pmatrix} \text{ for } j=1, \dots, d$$

$$\text{and prediction matrix: } Z_{(n \times d+1)} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ z_{1,1} & z_{1,2} & \dots & z_{1,d+1} \\ \vdots & \vdots & \ddots & \vdots \\ z_{n,1} & z_{n,2} & \dots & z_{n,d+1} \end{bmatrix}$$

thus, for polynomial regression  $Y \sim Z$ , then

$$\text{new } \hat{\beta}^* = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_d \end{pmatrix} = (Z^T Z)^{-1} Z^T Y$$

$$\text{where } Z \text{ is orthogonal s.t. } Z^T Z = \begin{pmatrix} \tilde{z}_{1,0}^2 & \tilde{z}_{1,1}^2 & \dots & 0 \\ \tilde{z}_{1,1}^2 & \tilde{z}_{1,2}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \tilde{z}_{1,d}^2 \end{pmatrix}$$

$$\Rightarrow \hat{\beta}_j^* = \frac{\tilde{z}_{1,j}^T Y}{\tilde{z}_{1,j}^T \tilde{z}_{1,j}} \text{ for } j=0, \dots, d$$

Practically, we can further normalize over  $\tilde{z}_j$ 's s.t.

$$\tilde{z}_{1,j}^T \tilde{z}_{1,j} = 1, j=0, \dots, d \text{ (T)}$$

use, for e.g.  $z_{1,0} = \frac{1}{\sqrt{n}}$  i.e.  $j=0, \dots, n$  (intercept)  
normalizing constant so eq holds

$$\Rightarrow \hat{\beta}_j^* = \tilde{z}_{1,j}^T Y \text{ for } j=1, \dots, d$$

Note: Practically, easier to safely choose  $d$  as highest-order significant polynomial term

(orthogonal  $\Rightarrow$  no correlation between different-order polynomial terms)

$\Rightarrow$  hierarchical structure simplified, i.e.

more unlikely to have contradictory results as seen w/ non-orthogonal preds)

### Polynomials in several predictors

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1^2 + \beta_4 X_2^2 + \beta_{12} X_1 X_2 + \dots + E$$

$d=2$  Interaction term

R: `lmCY ~ poly(x1, x2, deg=2)`

by default we're still fitting orthogonal polynomials

REMARKS: • can use F-test to select d for, e.g.

Model vs. Null:  $H_0: d=2$

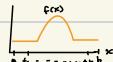
• For now, assume d is the same for all predictors we introduce ( $\geq 1$  pred.)

### Regression Splines (Nonparametric Statistics)

Disadvantage of polynomials: each data affects the fit globally

Remedy: B-spline

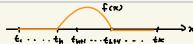
In general, suppose  $E(Y|X) = f(X)$  where



Cubic (B-spline) basis functions on interval  $[a, b]$  w/ pre-specified

Knots  $t_0, \dots, t_k$ :

- non-zero on interval defined by four successive knots
- and zero elsewhere  $\Rightarrow$  local influence property



where, e.g.:

$$\begin{cases} x_1 > t_2 & \Rightarrow f_1(x_1) = 0 \\ x_2 < t_1 \text{ or } x_2 > t_3 & \Rightarrow f_2(x_2) = 0 \\ \vdots & \vdots \\ x_4 > t_3 & \Rightarrow f_4(x_4) = 0 \end{cases}$$

s.t. different h gives different predictors  
where  $f_h$  takes cubic polynomial form that we skip

Y(h), f\_h

$$Y \sim f_1(x) + \dots + f_n(x) \quad \text{new predictors}$$

- cubic polynomial fit to each four successive knots
- smooth
- integrates to one

### Other Transformations

- Smoothing splines
- Generalized additive models (GAMs)
- CART, MARS, MABT, neural networks

#### Rule of thumb:

- For large datasets, complex models are better (w/ appropriate control of the number of parameters)
- For small datasets or high noise levels (e.g. social sciences), standard regression is more appropriate

### Ch: Factors (Qualitative Variables):

#### Material covered

##### Factors (Qualitative predict.)

- reg w/ one factor
- reg w/ one factor and common cont. preds
- main effects model

##### Multiple Factors

##### F-TEST

#### Outline

- Regression w/ one factor only
- Regression w/ one factor and other (continuous) predictors

#### Factors

- Qualitative (categorical) predictors; can have two/more levels

#### One-Factor Models

- Factor vars can be included in a MLE mean function using dummy variables

Suppose  $d=2$ : single dummy var:  $\begin{cases} 0 : \text{category 1} \\ 1 : \text{category 2} \end{cases}$

+ assignment of labels to values is generally arbitrary, and will not change the outcome of the analysis

+ can be alternatively defined w/ different set of values, e.g.

$\{1, -1, 2, 1, 2, 0, \dots\}$ , etc. (must have two distinct values)

Suppose  $d=3$ : j<sup>th</sup> dummy variable  $v_j$  for the factor,  $j=1, \dots, d$

has the i<sup>th</sup> value  $v_{ij}$  for  $i=1, \dots, n$ , i.e.

$$v_{ij} = \begin{cases} 1 & \text{if } \text{group}_j = j^{\text{th}} \text{ category of group} \\ 0 & \text{otherwise} \end{cases}$$

Ex (UN Data)observational study;  $n=199$ 

Response: lifeExp

Group: factor variable ( $d=3$ ):  
 $U_1 = \text{oecd}$   
 $U_2 = \text{other}$   
 $U_3 = \text{africa}$

$$\text{where } U_1 + U_2 + U_3 = 1$$

	Group	$U_1$	$U_2$	$U_3$
Afghanistan	other	0	1	0
Albania	other	0	0	1
Algeria	africa	0	1	0
Angola	africa	0	1	0
Anguilla	other	0	1	0
Argentina	other	0	1	0
Armenia	other	0	1	0
Aruba	other	0	1	0
Australia	oecd	1	0	0
Austria	oecd	1	0	0

Remark: if we add an intercept to the mean function, the resulting model is overparameterized.  
 $\Rightarrow U_1 + U_2 + U_3 = 1$ , a col. of 1's, and the column of 1's is the regressor that corresponds to int.

Account for this problem by dropping one of the dummy vars (e.g.  $U_3$ ):

$$E(Y) = \beta_0 + \beta_1 U_1 + \beta_2 U_2$$

$$\text{Since } U_1 + U_2 + U_3 = 1 \Rightarrow U_1 = 1$$

$$\Rightarrow E(Y|U_1=1, U_2=0, U_3=0) = \beta_0 + \beta_1(1) + \beta_2(0) = \beta_0 + \beta_1$$

$$\text{Similarly, } E(Y|U_1=0, U_2=1, U_3=0) = \beta_0 + \beta_1(0) + \beta_2(1) = (\beta_0 + \beta_2)$$

$$\text{and } E(Y|U_1=0, U_2=0, U_3=1) = \beta_0 + \beta_1(0) + \beta_2(1) = (\beta_0 + \beta_2)$$

Prop (2): Let Y-response (numerical), X-explanatory (categorical)

$$E(Y|X)$$

Table 5.1 Regression Summary for Model (5.4)

	Estimate	Std. Error	t-Value	Pr(> t )
(Intercept), $\beta_0$	82.4465	1.1279	73.09	0.0000
other, $\beta_1$	-7.1197	1.2709	-5.60	0.0000
africa, $\beta_2$	-22.6742	1.4200	-15.97	0.0000

$\hat{\sigma} = 6.2801$  with 196 df,  $R^2 = 0.6191$ .

3 groups are

$$\hat{E}(\text{lifeExp}|U_1=\text{oecd}) = \hat{\beta}_0 + \hat{\beta}_1(0) + \hat{\beta}_2(0) = 82.45$$

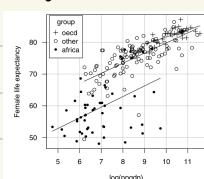
$$\hat{E}(\text{lifeExp}|U_1=\text{other}) = \hat{\beta}_0 + \hat{\beta}_1(1) + \hat{\beta}_2(0) = 82.45 - 7.12 \quad (5.5)$$

$$\hat{E}(\text{lifeExp}|U_1=\text{africa}) = \hat{\beta}_0 + \hat{\beta}_1(0) + \hat{\beta}_2(1) = 82.45 - 22.67$$

$\hat{E}(Y|U_1=\text{oecd}=1) = 82.45$ ; testing  $H_0: \beta_0 = 0$  not practically meaningful  
 $\hat{E}(Y|U_1=\text{other}=1) = 82.45 - 7.12$ ; testing for  $H_0: \beta_1 = 0$   
 $\hat{E}(Y|U_1=\text{africa}=1) = 82.45 - 22.67$ ; testing for  $H_0: \beta_2 = 0$  respectively

Remarks:  $df = n - (p+1)$ ; here,  $p=2 \Rightarrow df = 199 - (2+1) = 196$ 

Adding a continuous Predictor: suppose we add log(ppgdp) as continuous var.



can model w/ separate regressions for each factor (with

+ three lines w/ separate int. and slope

by group

Let group=j S.E.	int: oecd	slope
\$j=1\$: oecd	$\gamma_{01}$	$\gamma_{11}$
\$j=2\$: other	$\gamma_{02}$	$\gamma_{12}$
\$j=3\$: africa	$\gamma_{03}$	$\gamma_{13}$

$$\text{then } E(Y| \log(ppgdp) \geq x, \text{group}=j) = \gamma_{0j} + \gamma_{1j}x$$

where  $(\gamma_{0j}, \gamma_{1j})$  are the intercept and slope for level  $j = 1, \dots, d \Rightarrow \text{sd} = 6 \text{ parameters!}$ 

Instead, model is generally parameterized by

$$E(Y|\log(ppgdp) \geq x, \text{group}) = \beta_0 + \beta_{02}U_2 + \beta_{03}U_3 + \beta_1x + \beta_{12}U_2x + \beta_{13}U_3x \quad (*)$$

$$Y \sim \text{group} * \log(ppgdp) + \text{group} * \log(ppgdp)$$

|||

$$Y \sim \text{group} * \log(ppgdp) \quad \text{where } * \text{ expands to include all main effects / interactions}$$

$$\text{Oecd: } U_1=1, U_2=U_3=0 \Rightarrow E(Y|x, \text{oecd}) = \beta_0 + \beta_1x$$

$$\Leftrightarrow \gamma_{01} = \beta_0 + \beta_{02}, \gamma_{11} = \beta_1 + \beta_{12}$$

$$\text{Other: } U_2=1, U_1=U_3=0 \Rightarrow E(Y|x, \text{other}) = (\beta_0 + \beta_{02}) + (\beta_1 + \beta_{12})x$$

$$\Leftrightarrow \gamma_{02} = \beta_0 + \beta_{03}, \gamma_{12} = \beta_1 + \beta_{13}$$

$$\text{Africa: } U_3=1, U_1=U_2=0 \Rightarrow E(Y|x, \text{africa}) = (\beta_0 + \beta_{03}) + (\beta_1 + \beta_{13})x$$

$$\Leftrightarrow \gamma_{03} = \beta_0 + \beta_{02}, \gamma_{13} = \beta_1 + \beta_{12}$$

Remark: parameters  $(\beta_0, \beta_1)$  are the int and slope for the baseline level;  
 remaining  $\beta$ 's are the differences b/w the other levels & baseline

Table 5.3 Regression Summary for Model (5.7)

	Estimate	Std. Error	t-Value	Pr(> t )	
(Intercept), $\hat{\beta}_0$	59.2137	15.2203	3.89	0.0001	
2 <sup>nd</sup> order	other, $\hat{\beta}_2$	-11.1731	15.5948	-0.72	0.4746
	africa, $\hat{\beta}_3$	-22.9848	15.7838	-1.46	0.1470
	log(ppgdp), $\hat{\beta}_1$	1.5544	1.0165	1.53	0.1278
3 <sup>rd</sup> order	other: log(ppgdp), $\hat{\beta}_2$	0.6442	1.0520	0.61	0.5410
	africa: log(ppgdp), $\hat{\beta}_3$	0.7590	1.0941	0.69	0.4887
	$\hat{\sigma} = 5.1293$ with 193 df, $R^2 = 0.7498$ .				

Note: all terms are non-significant; likely due to correlation among terms

→ motivates removing higher-order (interaction) terms (i.e. main effects model)

Main Effects Model: set  $\beta_2 = \beta_3 = 0$   $E(Y|\log(ppgdp); x_1, x_2, x_3) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$   
 all 3 groups have same slope  $(\log group + \log(ppgdp))$

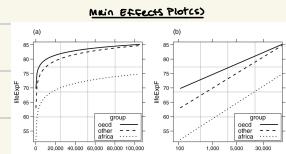


Figure 5.3: Effect plots for the interaction model (5.7) for the UN data: (a)  $y$  vs  $x_1$  on the log-linear axis; (b)  $y$  vs  $x_2$  in log-scale.

Remarks: - dropping int. terms from (8)

- allows each group to have its own intercept

- all groups have the same slope

+ difference b/wn the levels of the factor are the same for every fixed value of the cont. regressor

+ the effect of the cont. regressor is the same for all levels of the factor

	Estimate	Std. Error	t-value	Pr(> t )
(1 <sup>st</sup> order)	(Intercept), $\hat{\beta}_0$	49.5292	3.3996	14.57 0.0000
	other, $\hat{\beta}_2$	-1.5347	1.1737	-1.31 0.1926
	africa, $\hat{\beta}_3$	-12.1704	1.5574	-7.81 0.0000
	log(ppgdp), $\hat{\beta}_1$	2.2024	0.2190	10.06 0.0000

$\hat{\sigma} = 5.1798$  with 195 df,  $R^2 = 0.7422$ .

note: now significant

### Many Factors ( $> 2$ )

- increasing # of factors on a continuous Y in a mean function can add complexity; does not raise new fundamental issues

Table 5.5 The Wool Data

Variable	Definition
$y$	len
$x_1$	amp
$x_2$	load
$x_3$	log(cycles)

Note: there are  $3^3 = 27$  possible combinations of the 3 factors used.

Main Effects Mean  $E_y = Y \sim X_1 + X_2 + X_3$

where  $X_i$ : 2 dummy vars, 2=1, 2=3

$\Rightarrow (1+2+3) = 27$  params

Let  $X_1: U_{11}, U_{12}, U_{13}$ ,  $X_2: U_{21}, U_{22}, U_{23}$ ,  $X_3: U_{31}, U_{32}, U_{33}$

where  $U_{ij}$ : baseline/cmp. group,  $i=1, 2, 3$

$$E(Y|X_1, X_2, X_3) = \beta_0 + \beta_1 U_{11} + \beta_2 U_{21} + \beta_3 U_{31} + \beta_{12} U_{12} + \beta_{13} U_{13} + \beta_{23} U_{23}$$

$\Rightarrow E(Y|U_{11}=1, U_{12}=1, U_{13}=1) = \beta_0$  (all 3 factors set to low level)

$\Rightarrow E(Y|U_{11}=1, U_{12}=1, U_{13}=2) = \beta_0 + \beta_{12}$ , i.e.  $\beta_{12}$ : diff. b/wn  $U_{12}$  and  $U_{13}$  w.r.t.  $X_1$ , holding low level VALUES constant for  $X_1$  and  $X_3$

Full Second-order Mean  $E_y = Y \sim X_1 + X_2 + X_3 + X_1 X_2 + X_1 X_3 + X_2 X_3$

- adds all the two-factor interactions to the mean function  
 $(7+3 \cdot 4) = 19$  params

Full Third-order Mean  $E_y = Y \sim X_1 + X_2 + X_3 + X_1 X_2 + X_1 X_3 + X_2 X_3 + X_1 X_2 X_3$

- involves 3-factor interactions  
 $(10+8) = 18$  params

### F-tests (ch. 8)

- for any linear regression model, the RSS measures the amount of variation in response not explained by regressors

- if the null were false, then the RSS and  $\text{RSS}_{\text{null}}$  under null model

$\rightarrow$  basis of F-test: enough evidence against the null if difference  $(\text{RSS}_{\text{full}} - \text{RSS}_{\text{null}})$  is large enough

NH:  $E(Y|X_1, X_2, X_3) = \beta_0 + \beta_1 X_1$ ,  $\Rightarrow \beta_2 = 0$

Alt:  $E(Y|X_1, X_2, X_3) = \beta_0 + \beta_1 X_1 + \beta_2 X_2$

$$\text{anova } (\text{M}_{\text{HO}}, \text{M}_{\text{HOHA}}) \rightarrow \begin{array}{ccccccc} \text{Res.DF} & \text{RSS} & \text{DF} & \text{Sum of Square} & \text{F} & \text{Pr(>F)} \\ \text{dfHO} & \text{RSS}_{\text{HO}} & & & & & \\ \text{dfHA} & \text{RSS}_{\text{HA}} & \text{dfHO+dfHA} & \text{RSS}_{\text{HO+HA}} & & & \\ & & & & \text{F-stat} & \text{p-value} & \end{array}$$

where  $\text{M}_{\text{HO}} = \text{HO} + \text{HA}$

Note: t-test only appropriate for testing individual predictors (not, for e.g. interaction terms)

REMARKS: F-test requires nested structure, i.e. cannot test:

$$\text{H}_0: Y \sim X_1 + X_2 \text{ vs. } \text{H}_A: Y \sim X_1 + X_2 + \dots$$

## Ch. 10: Variable Selection

### Material covered

#### Testing-based approaches

- Backward elimination
- Forward selection
- Stepwise regression

#### Criterion-based approaches

- AIC and BIC
- Adjusted R<sup>2</sup>
- Mallows Cp

### Testing-Based Approaches

Idea: test significance of predictors and eliminate in some principled fashion

- Based on individual p-values
- multiple testing is not accounted for - but ranking is more important than the absolute size of p-values (not using FDR control) e.g. Bonferroni correction
- Different methods use different rules to add/delete predictors

#### Backward Elimination (assumes: n > p)

- Start w/ all predictors in the model
  - Remove the predictor w/ the highest p-value > a
  - REFIT model → return to 2
  - Stop when all p-values < a
- Note:  $a > 0.05$  may be better if prediction is the goal

#### Forward Selection ("Better" in practice for large p)

- Start w/ no pred. vars
- ↓ predictor not in model, check the p-value if added to model
- Add the one w/ the smallest p-value < a
- REFIT model and return to 2
- Stop when no new predictors can be added

$$\begin{array}{l} Y \sim X_1 + \text{Smallest p-value} \\ Y \sim X_1 \\ Y \sim X_1 + X_2 \\ \vdots \\ Y \sim X_1 + X_2 + \dots \\ Y \sim X_1 + X_2 + \dots + X_p \end{array}$$

Stepwise regression: combination of backward elimination and forward selection (allows to add variables back after they have been removed)

REMARKS: Greedy; may miss the optimal model,

e.g.  $Y \sim X_1 + \dots + X_p \Rightarrow 2^p$  possible (subset) models

$$\text{forward or backward} \leq \Pr(p=1) + (p-1) \cdot \frac{(p-1)p}{2}$$

at most

- Remember not to take p-values at face value (multiple testing)
- Variables not selected can still be correlated w/ the response, but they do not improve the fit enough to be included
- tend to pick smaller models than desirable for prediction purposes

### (Information) criterion-based selection procedure

Idea: choose the model that optimizes a criterion which balances goodness-of-fit (RSS) and model size (p)

- Implicitly uses either forward or's backward selection
- No p-values
- Some theoretical guarantees
- Different methods use different goodness-of-fit measures and different penalties for model size

#### AIC and BIC

$$\text{Akaike information criterion (AIC)}: = n \ln(\text{RSS}/n) + 2p\ln(n)$$

$$\Sigma \text{step}(\dots, k=2)$$

(Theoretically), suppose candidate model not  $X_{\text{true}}$ , then choose AIC

$$\text{Bayesian information criterion (BIC)}: = n \ln(\text{RSS}/n) + p \ln(n)$$

larger weight for p → often selects smaller models (minimizing)

$$\Sigma \text{step}(\dots, k \ln(n))$$

(theoretical) selection consistency property holds: if  $\underbrace{\text{candidate model}}_{X^{\text{true}}}$ , choose BIC

Remark: Practically, could get different model choices

Likelihood: A general def. of AIC/BIC:

For model C,  $AIC = -2\log(\text{max likelihood of model } C) + 2 \times \# \text{params in } C$

$BIC = -2\log(\text{max likelihood of model } C) + \log(n) \# \text{params in } C$

Linear regression models:

$$\text{Model (TRUE)} \left\{ \begin{array}{l} Y_i = X_i^T \beta + \epsilon_i, i=1, \dots, n \\ \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \text{ and } E(\epsilon_i) = N(0, \sigma^2 I_n) \end{array} \right.$$

$$\Leftrightarrow Y | X \sim N(X \beta, \sigma^2 I_n), i=1, \dots, n$$

The joint prob. density function (pdf) of  $Y = (Y_1, \dots, Y_n)^T$  given  $X$ :

$$\begin{aligned} p(Y_1, \dots, Y_n | X; \beta, \sigma^2) &= \prod_{i=1}^n p(Y_i | X_i; \beta, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp \left\{ -\frac{(Y_i - X_i^T \beta)^2}{2\sigma^2} \right\} = L(\beta, \sigma^2 | X, Y) : \text{function of } \beta, \sigma^2 \text{ given Data } (X, Y) \end{aligned}$$

The estimates maximizing the likelihood function  $L(\beta, \sigma^2 | X, Y)$  or, equivalently,

$$\log(\text{likelihood}) = \log L(\beta, \sigma^2 | X, Y) = \sum_{i=1}^n \left[ \log \left( \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \right) - \frac{(Y_i - X_i^T \beta)^2}{2\sigma^2} \right] = \frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - X_i^T \beta)^2$$

are called the maximum likelihood estimators (MLE):  $(\hat{\beta}_{MLE}, \hat{\sigma}_{MLE}^2)$

It can be shown that  $\text{OLS} \hat{\beta} = \hat{\beta}_{MLE}$ :

$$\text{to get } \hat{\beta}_{MLE}, \frac{\partial \log(L(\beta, \sigma^2 | X, Y))}{\partial \beta} = \frac{\partial}{\partial \beta} \sum_{i=1}^n (Y_i - X_i^T \beta)^2 = 0 \Rightarrow \hat{\beta}_{MLE} = \hat{\beta}_{OLS}$$

$$\hat{\sigma}_{MLE}^2 = \frac{\partial \log L}{\partial \sigma^2} = 0 \Rightarrow \frac{1}{2} + \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - X_i^T \beta)^2 = 0$$

$$\Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - X_i^T \hat{\beta}_{MLE})^2 = \frac{RSS}{n}$$

$\neq \hat{\sigma}^2 = \frac{RSS}{n-p+1}$ ; but,  $\hat{\sigma}^2 - \hat{\sigma}_{MLE}^2 \xrightarrow{n \rightarrow \infty} 0$  w.r.t.  $|p \text{ fixed}|$

$$E(\hat{\sigma}^2) = \sigma^2, \text{ i.e. unbiased and bias: } \hat{\sigma}_{MLE}^2 \xrightarrow{n \rightarrow \infty} 0.$$

Back to AIC/BIC:

$$\begin{aligned} \text{For linear model } C, \log(\text{max likelihood model } C) &= \frac{n}{2} \log((2\pi\hat{\sigma}_{MLE}^2)^{-1}) - \frac{1}{2\hat{\sigma}_{MLE}^2} \sum_{i=1}^n (Y_i - X_i^T \hat{\beta}_{MLE})^2 \\ &= -\frac{n}{2} \log \left( \frac{RSS_C}{n} \right) - \frac{n}{2} + \text{const.} \quad |_{n \rightarrow \infty} \\ &\approx n \log \left( \frac{RSS_C}{n} \right) + \text{const} \end{aligned}$$

Adjusted R<sup>2</sup>

$$\text{Recall: } R^2 = 1 - \frac{RSS}{TS}$$

$$\begin{aligned} \text{adjusted } R^2 &:= R^2_C = 1 - \frac{RSS_C / (n-p+1)}{TS_C / (n-p+1)} \\ &= 1 - \left( \frac{RSS_C}{n-p+1} \right) (1-R^2) \end{aligned}$$

• adding a predictor will not necessarily increase R<sup>2</sup>

• Maximizing R<sup>2</sup> is equivalent to minimizing SEE  $\hat{\sigma}$

Mallows' Cp

$$C_p := \frac{RSS_P}{\hat{\sigma}^2} + 2(p+1-n)$$

•  $C_p \leq p+1$  indicates good fit

• Cp estimates the (test data) MSE:  $f \in E(Y_i - \hat{Y}_i)^2$

[suppose we have another set of data w/ same  $X^*$ ]; then,  $MSE = \frac{1}{n} \sum_i (Y_i^* - \hat{Y}_i^*)^2$   
as original data but  $Y^*$ 's new responses

Remark: Information criteria can be used together with:  
 - forward selection  
 - backward selection  
 - best subset selection

### Variable Selection Summary

• Sensitive to outliers

• Generally, criterion-based methods are preferred

• It may happen that several models provide very similar fit

• If models w/ similar fit lead to very different conclusions, the data are ambiguous

• If conclusions are similar, choose a simpler model and/or predictors that are easier to measure

ch11: shrinkage methods:

- outcome
- Ridge Regression
- LASSO
- (skipping PLS and PCA)

### Ridge Regression

Penalizing the square of the coefficients:

$$\min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (1) \text{; Lagrangian function for } \lambda \text{- see below}$$

$\underbrace{\text{RSS}(\beta)}$

coefficients  $\hat{\beta}_{\text{ridge}}$  are shrunken towards zero

$\lambda > 0$  is tuning parameter (estimated)

$\beta_0$  controls the amount of shrinkage:  $\begin{cases} \beta_0: \text{OLS estimation, i.e. } \hat{\beta}_{\text{ridge}} \rightarrow \beta_0 \text{ OLS} \\ \lambda > 0: \hat{\beta}_{\text{ridge}} = 0 \text{ for } j=1, \dots, p \\ \text{(since otherwise } \lambda \sum_j \beta_j^2 \rightarrow \infty) \end{cases}$

note: if  $\beta_j = 0$  for  $j=1, \dots, p$ , the (1) becomes  $\min_{\beta_0} \sum_i (y_i - \beta_0)^2 \Rightarrow \hat{\beta}_0^{\text{ridge}} = \bar{y}$

$$\Leftrightarrow \hat{\beta}_{\text{ridge}} \rightarrow \bar{y} \text{ as } \lambda \rightarrow 0$$

Equivalent Formulation:  $\min_{\beta} \sum_i (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq S \quad (2)$

where  $S$  (shrinkage parameter) s.t.:  $\begin{cases} S > 0: \hat{\beta}_{\text{ridge}} \neq 0 \text{ for } j=1, \dots, p \text{ and } \hat{\beta}_0^{\text{ridge}} \rightarrow \bar{y} \\ (S=0) \text{ or } (\lambda=0): \hat{\beta}_{\text{ridge}} = \beta_{\text{OLS}} \end{cases}$

why? when there are many highly correlated variables (collinearity)

$\hat{\beta}_{\text{OLS}}$  may have large coefficient on one variable and a similarly large negative coefficient on its correlated variable (collinearity)

in ridge regression, the size constraint tries to avoid this phenomenon

address (fit model) under high-dimensional problem ( $p > n$ ) when OLS fails ( $X^T X$  not invertible)

remark: ridge estimate is not equivalent under scaling of the predictors

In practice, standardize the predictors first:  $x_{ij} \leftarrow \frac{x_{ij} - \bar{x}_i}{\text{SD}(x_{ij})} \text{ for } j=1, \dots, p$

similarly, can center (or standardize)  $y_i - y_i - \bar{y}$  or  $\leftarrow \frac{y_i - \bar{y}}{\text{SD}(y)}$  for  $i=1, \dots, n$

Ridge Estimation: Suppose  $X$ -standardized and  $y$ -centered (standardized).

$$\hat{\beta}_{\text{ridge}} = 0 \text{ since } \partial \text{objective (1)} / \partial \beta_0 = 0$$

$$\Rightarrow -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip}) = 0$$

$$\Rightarrow (\sum_i y_i) - n \beta_0 - \beta_1 \sum_i x_{i1} - \dots - \beta_p \sum_i x_{ip} = 0$$

use  $y$ -centered,  $x$ -standardized  $\Rightarrow \sum_i y_i, \sum_i x_{ij} = 0$

$$\Rightarrow \hat{\beta}_{\text{ridge}} = 0.$$

2) the ridge problem is equivalent to the following w/o intercept:

$$\hat{\beta}_{\text{ridge}} = \underset{\beta_{\text{real}}} {\text{argmin}} \sum_{j=1}^p \beta_j^2 \text{ where } \beta = (\beta_1, \dots, \beta_p)^T \quad \begin{array}{l} \text{no intercept} \\ x_{0j} = 1 \text{ for all } j \end{array}$$

$$= \underset{\beta_{\text{real}}} {\text{argmin}} (X^T X \beta)^2 + \lambda \beta^T \beta \text{ where } y = (y_1, \dots, y_n)^T$$

$$X = (x_{ij})_{n \times p} = \begin{pmatrix} 1 & \dots & 1 \\ x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix}$$

solving as done w/ OLS.

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

Remarks:  $\cdot X$ -standardized,  $y$ -centered (standardized)

$\cdot \beta_{\text{OLS}}$  and  $\hat{\beta}_{\text{ridge}}$  will be different as above

$\cdot \hat{\beta}_{\text{ridge}}$  is linear in  $y$  (as seen in OLS)

$\cdot \hat{\beta}_{\text{ridge}}$  is biased (less unbiased);  $E(\hat{\beta}_{\text{ridge}}) \neq \beta$  (unbiased) for  $\lambda > 0$

$\cdot$  Even if  $X$  is not full rank,  $(X^T X + \lambda I)$  is invertible

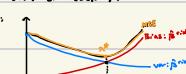
$\cdot \hat{\beta}_{\text{ridge}}$  has smaller variance than OLS; thus, may have smaller MSE

MSE: for a new obs.  $(x^*, y^*)$ ,

$$E(y^* - \hat{y}^*)^2 = \text{Var}(\hat{y}^*) + [E(\hat{y}^*)]^2 + \text{const.} \xrightarrow{\text{ignores}}$$

$$\text{use } \hat{y}^* = X^T \hat{\beta}_{\text{ridge}} \text{ and } \text{Var}(\hat{y}^*) = E(\hat{y}^*)^2 - E(\hat{y}^*)^2 = (X^T \hat{\beta}_{\text{OLS}} - \beta)^T (X^T \hat{\beta}_{\text{OLS}} - \beta)$$

ignores uncorr. bias term(s):



Shrinkage in Ridge: Suppose orthonormal design ( $X^T X = I$ ). Then,  $\hat{\beta}^{OLS} = X^T Y$  and

$$\begin{aligned} RSS(\beta) &= (Y - X\beta)^T (Y - X\beta) = Y^T Y - 2Y^T X\beta + X^T X\beta \\ &\quad \text{constant} \quad \underbrace{I}_{X^T X} \\ &= \|Y\|^2 - 2(Y^T X)\beta + \text{const} \\ &= (\beta - \hat{\beta}^{OLS})^T (\beta - \hat{\beta}^{OLS}) + \text{const} \\ &= \text{const.} + \sum_{j=1}^p (\beta_j - \hat{\beta}_j^{OLS})^2 \end{aligned}$$

$$\Rightarrow (1) w/ an orthogonal design \equiv \sum_{j=1}^p (\beta_j - \hat{\beta}_j^{OLS})^2 + 2 \sum_{j=1}^p \hat{\beta}_j^{OLS}$$

equivalent to component-wise minimization:

$$\min_{\beta_j} (\beta_j - \hat{\beta}_j^{OLS})^2 + \lambda \beta_j^2 \rightarrow \frac{\partial \text{RSS}}{\partial \beta_j} = 0 \Rightarrow \hat{\beta}_j^{\text{Ridge}} = \frac{1}{1+\lambda} \hat{\beta}_j^{OLS}$$

Remarks (Shrinkage in Ridge): - shrink the estimate towards zero by a positive constant less than 1.

$$\text{Var}(\hat{\beta}_j^{\text{Ridge}}) = \frac{1}{(1+\lambda)} \text{Var}(\hat{\beta}_j^{OLS}) ; \text{Bias}(\hat{\beta}_j^{\text{Ridge}}) = \beta_j - \frac{\lambda}{1+\lambda} \beta_j$$

E(BIAS):  $E(\hat{\beta}_j^{\text{Ridge}}) - \beta_j$

$$= E\left(\frac{\hat{\beta}_j^{OLS}}{1+\lambda}\right) - \beta_j$$

+ 2.7, shrinkage P, bias P, variance L

$$+ 2.4, \text{ shrinkage L, bias L, variance L}$$

### Model Assessment

- objectives: (1) choose a value of a tuning parameter ( $\lambda$ ) for a technique
- (2) estimate the prediction performance of a given model
- run procedure on an independent test set (if available)
- one should use different test data for (1) and (2); a validation set for (1)
- a test set for (2)

### Cross-validation

- often insufficient data to create a separate validation or test set; setting some data aside for validation is possible, but affects the accuracy of training estimates



K-Fold CV : (1) divide the data into K disjoint subsets

(2) use subsets 1,...,K as training data and subset K as validation data  
compute the MSE on subset K, fitting over complete models (different  $\lambda$ 's)

(3) repeat for each subset

(4) average the result

Data	
y	Sub 1
x <sub>1</sub>	Sub 2
⋮	⋮
x <sub>p</sub>	Sub K

### Least Absolute Shrinkage and Selection Operator (LASSO)

$$\min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + 2 \sum_{j=1}^p |\beta_j| \quad \text{Gives } \min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 \text{ subject to } \sum_{j=1}^p |\beta_j| \leq S \quad \text{where } S \geq 0 \text{ (tuning parameter) s.t.: } \begin{cases} S > 0, \hat{\beta}_j^{\text{Lasso}} = 0 \text{ for } j=1, \dots, p \text{ and } \hat{\beta}_0^{\text{Lasso}} = \bar{y} \\ S = 0 \text{ (} \lambda = 0 \text{)}: \hat{\beta}^{\text{Lasso}} = \hat{\beta}^{\text{OLS}} \end{cases}$$

• shrinkage:  $\hat{\beta}_j^{\text{Lasso}} = 0$  for  $|\beta_j^{\text{OLS}}| \geq S$ ,  $\hat{\beta}_j^{\text{Lasso}} = \beta_j^{\text{OLS}}$  for  $|\beta_j^{\text{OLS}}| < S$

and: same as Stein's Ridge regression

• sparsity: some fitted coefficients are exactly zero

similar to OLS

+ motivates continuous variable selection

### Soft Thresholding

Practically, we assume X-standardized, Y-centered (standardized)

$\hat{\beta}_j^{\text{Lasso}} = 0$  and we don't consider the intercept term

when X is orthonormal ( $X^T X = I$ ), we can minimize over  $\beta$  components:

$$\min_{\beta} |\beta_j - \hat{\beta}_j^{\text{OLS}}|^2 + \lambda |\beta_j| \quad \text{Solution: } \hat{\beta}_j^{\text{Lasso}} = \begin{cases} \hat{\beta}_j^{\text{OLS}} \pm \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{OLS}} \pm \frac{\lambda}{2} \neq 0 \\ 0 & \text{if } |\hat{\beta}_j^{\text{OLS}}| \pm \frac{\lambda}{2} = \text{Sign}(\hat{\beta}_j^{\text{OLS}}) \cdot (\hat{\beta}_j^{\text{OLS}} \pm \frac{\lambda}{2}) \\ \hat{\beta}_j^{\text{OLS}} \pm \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{OLS}} \pm \frac{\lambda}{2} = 0 \end{cases}$$

• shrinks large coefficients by a constant ( $\lambda/2$ )

• truncates small coefficients to zero

Remarks: other popularly implemented approaches: group Lasso (adding cat. vars as one group)

elastic net (combine L1/L2 norm)

### Summary

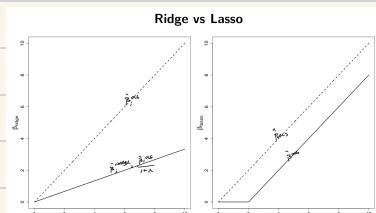
• main reason to use shrinkage: too many predictors or collinearity

• interpretation is usually lost

• range is still a linear model in w/ standardized (centered) pres./response original predictors but no inference

• prediction is usually improved by shrinkage

• all require selecting a tuning parameter



## Ch: Binomial Data

### Outline

- Binomial data
- Generalized Linear Models (GLM) for Binomial data
- Inference for GLM for binomial data
- Odds Ratio
- Overdispersion

### Review: the Binomial Distribution

$n$ : independent trials;  $i = 1, \dots, n$

$$\text{Suppose: } \begin{cases} P(Z_i=1) = p & (\text{success}) \\ P(Z_i=0) = 1-p & (\text{failure}) \end{cases}$$

Binomial RV:  $Y = \sum_{i=1}^n Z_i$  is the number of successes out of  $n$  iid trials.

$$\text{Prob Dist: } P(Y=k) = \binom{n}{k} p^k (1-p)^{n-k}, k=0, 1, \dots, n$$

$$\cdot E(Y) = np, V(Y) = np(1-p); E(Z_i) = p, V(Z_i) = p(1-p) \text{ for } i=1, \dots, n$$

$$\cdot \text{as } n \rightarrow \infty, \text{ Binomial} \xrightarrow{\text{CLT}} \text{Normal} \sim N(\mu, \sigma^2) \text{ where } \mu = np, \sigma^2 = np(1-p)$$

$$\cdot \text{Sample proportion estimate of } p: \hat{p} = \frac{Y}{n}$$

### Binomial Data: $i=1, \dots, N$

• response  $y_i$ : number of successes out of  $n_i$  independent trials w/ prob. of success  $p_i$

$$\text{where } y_i \sim \text{Binomial}(n_i, p_i), E(y_i) = n_i p_i$$

$x = (x_1, x_2, \dots, x_p)$  s predictors (quantitative, factors, or both)

assume: For all trials contributing to one response  $y_i$ , the predictors  $x_i$  have the same value (covariate class)

Goal: Model the relationship between  $y$  and  $x_1, \dots, x_p$  via modeling the relationship between  $p_i$  and  $x_1, \dots, x_p$

$$\text{where } p_i = \frac{E(y_i)}{n_i}$$

### Binomial Regression

$$\text{assume: } \begin{cases} y_i \sim \text{Binomial}(n_i, p_i) \\ y_i \text{ are independent} \end{cases}$$

want to construct a linear predictor, i.e.

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

remark: cannot use  $y_i = p_i$ ; need  $\text{prob}(y_i)$

Idea: use a link (transformation) function

$$y_i = g(p_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}, p_i \in \mathbb{R}$$

Binomial link functions:

$$\text{logit: } \eta = \log\left(\frac{p}{1-p}\right)$$

$$\Rightarrow \text{logistic regression s.t. } \log\left(\frac{p_i}{1-p_i}\right) = \eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

where odds:  $\frac{p}{1-p}$  conducive to interpretation

$$\text{Probit: } \eta = \Phi^{-1}(p), \Phi: \text{cdf of } N(0,1)$$

$$\Rightarrow \Phi^{-1}(p_i) = \eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

consider binary response, i.e.  $n_i=1$ . Then, using probit link

is equivalent to  $y_i = \begin{cases} 1 & \text{if } \eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} > 0 \text{ where } \eta_i \sim N(0,1) \\ 0 & \text{otherwise} \end{cases}$

$$\begin{aligned} \Rightarrow P_i = P(Y_i=1) &= P(X_i \beta + \epsilon_i > 0) \\ &= P(-\beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip} < 0) \\ &= \Phi(X_i \beta), \Phi: \text{cdf of } N(0,1) \end{aligned}$$

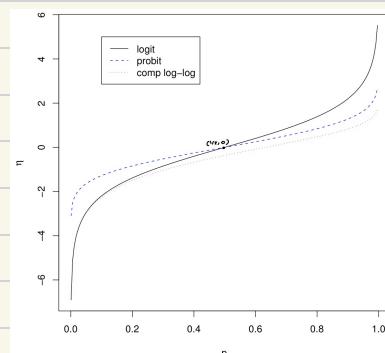
$$\Rightarrow \Phi^{-1}(p_i) = X_i \beta$$

\* incorporates hidden variable, but permits continuous response linear regression

$$\text{complementary log-log: } \eta = \log(-\log(1-p))$$

(less common in practice)

note: all transform  $P(Y_i=1) \rightarrow \eta \in (-\infty, \infty)$



Estimating parameters: assume logistic model:

$$\text{for } i=1, \dots, N, Y_i \sim \text{Binomial}(n_i; p_i) \text{ with } \log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_i \text{ so } p_i = \frac{\exp(\beta_0 + \beta_1 x_i)}{1 + \exp(\beta_0 + \beta_1 x_i)}$$

Maximum likelihood approach: find parameters ( $\beta_0, \beta_1, p_i$ ) that maximize the likelihood of the data, i.e.

$$\begin{aligned} & \prod_{i=1}^N P(Y_i=y_i) \text{ where } Y_i \sim \text{Binomial}(n_i; p_i) \\ & = \prod_{i=1}^N \left( \frac{n_i!}{y_i!(n_i-y_i)!} \right) p_i^{y_i} (1-p_i)^{n_i-y_i} \\ & L(p_1, \dots, p_N; y) = \sum_{i=1}^N [\log\left(\frac{p_i}{1-p_i}\right) + y_i \log(p_i) + (n_i - y_i) \log(1-p_i)] \\ & \text{and substitute } p_i \text{ via corresponding sigmoid function: } \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \\ & \Rightarrow L(\beta_0, \beta_1; y) \text{ where } \beta = \text{only remaining unknown quantity} \end{aligned}$$

Need to maximize

$$L(\beta) = \sum_{i=1}^N [\beta_0 (n_i \beta_1 x_i) - n_i \log(1 + \exp(\beta_0 + \beta_1 x_i))] + \text{const} \text{ (does not depend on } \beta)$$

obtain  $\hat{\beta}_{MLE}$  via optimization (numerical) algorithm; no closed form solution

Note: if data is perfectly separable,  $\hat{\beta}_{MLE}$  is not identifiable;  
consider a shrinkage procedure (e.g. ridge regression) to introduce a penalty/constraint for shrinking

## Inference

$$\begin{cases} H_0: \beta_j = 0 & \text{for } j=1, \dots, p \text{ predictors} \\ H_a: \beta_j \neq 0 \end{cases}$$

where  $\hat{\beta}_j := \frac{\hat{\beta}_j}{\text{se}(\hat{\beta}_j)}$

Note: similar to t-test in linear regression

## Likelihood ratio test

- two nested models
- L is the larger model w.r.t. parameters and likelihood  $L_L | \hat{\beta}_{MLE}$ ; model L
- S is the smaller model w.r.t. parameters and likelihood  $L_S | \hat{\beta}_{MLE}$ ; model S
- the likelihood ratio test statistic:  $= 2 \log \frac{L_L}{L_S} = 2[\log L_L - \log L_S] \stackrel{H_0}{\approx} \chi^2_{d.f.}$

Note: similar to (Covariates) F-test in OLS; measures difference in RSS under  $H_0$  vs.  $H_a$

$$\text{For logistic regression, LRT} = 2 \sum_{i=1}^N [Y_i \log \frac{\hat{p}_i^L}{\hat{p}_i^S} + (n_i - Y_i) \log \frac{1 - \hat{p}_i^L}{1 - \hat{p}_i^S}]$$

$$\text{where } \hat{p}_i^L = \frac{\exp(\beta_0^L + \beta_1^L x_i)}{1 + \exp(\beta_0^L + \beta_1^L x_i)},$$

$$\hat{p}_i^S = \frac{n_i \hat{p}_i^L}{n_i \hat{p}_i^L + n_i - \hat{p}_i^L},$$

$$\text{and } \frac{(1 - \hat{p}_i^L)}{(1 - \hat{p}_i^L)} = \frac{n_i - \hat{p}_i^L}{n_i - \hat{p}_i^L}$$

$$\begin{aligned} \text{Because: } \frac{\hat{p}_i^L}{1 - \hat{p}_i^L} & \xrightarrow{n \rightarrow \infty} \chi^2_{d.f.} \\ \frac{\hat{p}_i^L}{n_i \hat{p}_i^L} & \xrightarrow{n \rightarrow \infty} \text{constant by LLN} \end{aligned}$$

## Deviance

analog to RSS in OLS setting

take larger model to be the saturated model:  $= Y_i \stackrel{i.i.d.}{\sim} \text{Binomial}(n_i; p_i)$ ,  $i=1, \dots, N$ ;

N parameters:  $p_1, \dots, p_N$  to fit each datapoint perfectly,  
with fitted values  $\hat{p}_i = Y_i / n_i$

in this case, the (test statistic) is called the deviance of S

$$D := 2 \sum_{i=1}^N [Y_i \log \frac{Y_i}{n_i} + (n_i - Y_i) \log \frac{n_i - Y_i}{n_i}]$$

model S (candidate/consider) model compared to saturated model

where  $y_i = 1, \hat{p}_i, p_i$  are fixed prob's from S

Prop: If  $Y_i$ 's are binary binomials, independent,  $n_i$  are large

then  $D \stackrel{H_0}{\approx} \chi^2_{d.f.}$

$$\text{where: } \begin{cases} \text{H}_0: S \text{ model} \\ \text{H}_a: \text{saturated model} \end{cases}$$

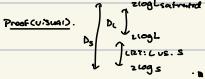
Remarks (uses of deviance): D tests goodness-of-fit:

$$p\text{-value} = P(\chi^2_{n-s} > D) ; \text{small } D \Rightarrow (\text{1}) \text{ good fit} \\ \text{large } D \Rightarrow (\text{2}) \text{poor fit}$$

• compare two nested models, e.g. null (no predictors) and current model;  
in this case, use

$$2 \log \frac{L_c}{L_0} = D_c - D_0 \approx \chi^2_{m+2-m-1} = \chi^2_{s-s}$$

where p-value:  $P(\chi^2_{s-s} > D_c - D_0)$



• (If  $n_i < 1$  binary data), deviance cannot be used (for goodness-of-fit);  
but GOF can still be used to compare model vs. L (S.L. fixed)

• From R output: null deviance (from model w/ only intercept)  
residual deviance (from fitted model)

Note: AIC is a general metric to compare all (including non-nested) models

#### Other measures of goodness-of-fit

• help analyze contingency table  
(categorical data) analysis, e.g.

	$y_{ij}$	$m_{ij}$
	.	.
	.	.
$\sum_j y_{ij} = Y_{i\cdot}$	$\sum_j m_{ij} = M_{i\cdot}$	
$\sum_i y_{ij} = y_{\cdot j}$	$\sum_i m_{ij} = M_{\cdot j}$	
$\sum_{ij} y_{ij} = N$	$\sum_{ij} m_{ij} = M$	

Here,  $\chi^2 = \sum_{i,j} \left[ \frac{(y_{ij} - m_{ij})^2}{m_{ij}} + \frac{(M_{i\cdot} - Y_{i\cdot})(M_{\cdot j} - y_{\cdot j})}{M} \right]$

$$\text{Success} = \frac{\sum_{i,j} y_{ij}}{N} = \frac{N}{M}$$

$$\text{Failure} = \frac{\sum_{i,j} m_{ij}}{M} = \frac{M}{N}$$

defined  $\chi^2$  goodness-of-fit statistic, i.e.

$$\text{Pearson's } \chi^2 := \chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i} \quad \begin{array}{l} \text{test: H}_0: \text{fitted model} \\ \text{vs. H}_a: \text{saturated model (least restriction)} \end{array}$$

use  $\begin{cases} O_i: \text{observed count in each "bin" (e.g. } y_{ij}\text{)} \\ E_i: \text{expected count under tested model (e.g. } \hat{y}_{ij}\text{)} \end{cases}$

For binomial data, add successes and failures to get

$$\chi^2 = \sum_{i=1}^n \frac{(y_{i\cdot} - n_i p_i)^2}{n_i p_i (1-p_i)} \Leftrightarrow \text{Pearson residuals: } r_i = \sqrt{\frac{(y_{i\cdot} - n_i p_i)^2}{n_i p_i (1-p_i)}} \\ \Rightarrow \chi^2 = \sum_{i=1}^n (r_i)^2$$

Remark: typically  $\chi^2$  is deviance (D) and is practically used the same way, i.e.

$$\chi^2 \stackrel{n}{\sim} \chi^2_{n-s} = D$$

#### Fidence intervals for Parameters

Asymptotically,  $\hat{\beta}$  is normal, i.e.  $\hat{\beta} - \beta_{\text{true}} \approx N(0, V)$

$\Rightarrow$  can use  $\beta$ -intervals; 95% CI:  $[\hat{\beta} \pm 1.96 \text{SE}(\hat{\beta})]$

Profile likelihood CI's are more accurate (based on considering the likelihood of the parameter w/ all others fixed)

#### Confidence Intervals for Predictions

- Predict probability of success (prob) :=  $\frac{\exp(x_i^T \hat{\beta})}{1 + \exp(x_i^T \hat{\beta})}$  for a particular  $x_i$
- no distinction b/w future observations & mean response
- Based on asymptotic normality of  $\hat{\beta}$  and  $x_i^T \hat{\beta}$
- Extrapolation will give unreliable predictions (as always)

#### Interpreting odds

$$\text{odds} := \frac{p}{1-p}$$

logistic regression (log-linear) models log odds, i.e.

$$\text{log odds} := \log \left( \frac{p}{1-p} \right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

interpretation: a unit increase in  $x_i$ , w/ all other predictors held fixed leads to an increase of  $\beta_i$  in log-odds; equivalently, odds multiplied by  $e^{\beta_i}$  (where  $e$  is reference value)

• no such interpretation available for other link functions

#### What does a large deviance indicate?

- violation of model assumptions (outliers, non-linearity, model structure)
- sparse data (small  $n_i$ )
- overdispersion

## Overdispersion

The binomial model links the mean and the variance:

$$\text{Var}(Y_i) = n_i p_i (1-p_i); \quad p_i \text{ parameter determines E(Y}_i\text{) and Var(Y}_i\text{)}$$

(not the case for normal data)

Overdispersion := observed  $\text{Var}(Y_i)$  is larger than the model postulates ( $\frac{\mu}{n}$ )

Common causes:

- | trials are not independent
- | prob. of success is not constant

Note: underdispersion is also possible but rare in practice

### Estimating overdispersion

Introduce an additional  $\boxed{\text{dispersion parameter} := \phi = \sigma^2}$

$$\text{s.t. } \text{Var}(Y_i) = \sigma^2 n_i p_i (1-p_i)$$

can estimate  $\sigma^2$  (as in linear regression) as

$$\hat{\sigma}^2 = \frac{\chi^2}{n - (p + 1)}$$

Remarks:

- does not affect  $\hat{p}$

- all standard errors must be multiplied by  $\sqrt{\phi}$

- Deviance can no longer be used to compare models (Binomial assumption no longer strictly met)

Instead use an

$$\boxed{\text{approximate F-test: } \frac{(D_0 - D_1)}{\hat{\sigma}^2} / (df_0 - df_1, n - p + 1)}$$

where  $D_0, D_1, df_0, df_1, n - p + 1$

- Goodness-of-fit cannot be tested (same reason)

- Estimating overdispersion is only reasonable when  $n_i$ 's are roughly equal

### Summary

- w/ link functions, binomial data can be modeled easily
- approximate inference available for testing models and parameter values
- logit has advantages in interpretation

Warnings: the estimation algorithm may not converge

· w/ small  $n_i$ , the  $\chi^2$  approx. is poor

overdispersion can be accounted for - but binomial assumption is sacrificed

## Ch: count regression

### Outline

- Review of the Poisson Dist.
- Poisson regression
- Inference via deviance
- Overdispersion
- Ex: Overcount data

Review: the Poisson Dist.

RV  $Y$  takes values  $0, 1, 2, \dots$

Poisson Dist. has one parameter:  $m \geq 0$

$$\boxed{\text{pdf: } P(Y=y) = \frac{e^{-m} m^y}{y!}} \quad E(Y) = \text{Var}(Y) = m$$

can be used to approximate Binomial( $n, p$ ) if  $\frac{n}{p} \gg m$  then  $np \gg m \gg 1$

If  $m$  is large,  $Y$  is approx. normal, i.e.  $\frac{Y-m}{\sqrt{m}} \rightarrow N(0, 1)$  as  $m \rightarrow \infty$

If  $Y_i \sim \text{Poisson}(m_i)$ ,  $Y_i$  - independent then  $\sum Y_i \sim \text{Poisson}(\sum m_i)$  where  $m = \sum m_i$

### Modelling count data

Response  $y_i$  is a count's count dist. to assume?

- If count is bounded above, binomial might be more appropriate
- If counts are large, normal approx. applies and regular linear regression may be used
- If count arises as number of 'failures' w/in a given number of 'successes', then negative binomial is appropriate

