

**Transformations**

Box-Cox works well for monotonic/simple relationships

**Box-Cox:**  $P(X) = \begin{cases} \frac{X-1}{P-1}, P \neq 0 \\ 1, P=0 \end{cases}$

**L'Hopital's Rule:**  $\lim_{P \rightarrow 0} \frac{X-1}{P-1} = \lim_{P \rightarrow 0} \frac{1}{P} = \infty$

**ex)**  $\text{Box-Cox}(P, N) = \frac{X-1}{P} \rightarrow \frac{(X-\lambda)^P-1}{P}$

- what is the reason for dividing by  $P$ ?  $P \rightarrow 0$ ,  $\lim_{P \rightarrow 0} \frac{X^P-1}{P} = \ln(X)$  so Box-Cox is continuous in  $P$ .

-  $X$  shift by 1 ensures that all transformations go through (0,1).

[Metric for symmetry]:  $S = \frac{\text{upper quartile} - \text{median}}{\text{median} - \text{lower quartile}}$

$S > 1$ : skewed right

for monotonically increasing transformations,  $\text{Var}(f(X)) = f(\text{Var}(X))$   
 $\text{median}(f(X)) = f(\text{median})$  for all quartiles.

only when

- all  $X$  values are positive  $\rightarrow \sqrt{X}, \log, \text{inverse}$  are undefined for negative or zero values
- power transformations are not monotone (e.g.  $X^2$ )

**solution:** If data domain includes 0 or is negative or data are all large or close together, then shift values by some  $\lambda$

2) ratio of the largest data to smallest data is sufficiently large. If ratio  $\approx 1$ , then the power transformation is roughly linear, and ineffective at bending the data. Shifting by  $\lambda$  solves the problem.

	$\log_{10} X$
2011	3.30341
2012	0.00022
2013	3.30363
2014	3.30384
2015	0.00022

	$\log_{10} X - 2010$
2011	0
2012	0.301
2013	0.176
2014	0.477
2015	0.125

	$\log_{10} X$
2011	0.097
2012	0.602
2013	0.602
2014	0.602
2015	0.609

Power transformations preserve the order of the data only when all values are positive and are effective only when the ratio of the largest to the smallest data values is itself large. When these conditions do not hold, we can impose them by adding a positive or negative start to all the data values.

Descending the ladder of powers (e.g., to  $\log X$ ) tends to correct a positive skew; ascending the ladder of powers (e.g., to  $X^2$ ) tends to correct a negative skew.

**positive skew**  $\log_{10} X$  **log transformation pulls in the right tail.**

	$\log_{10} X$
1	0
10	1
100	2
1000	3

**Transforming Nonlinearity**

Simple monotone nonlinearity can often be corrected by a power transformation of  $X$ , of  $Y$ , or of both variables. Mosteller and Tukey's bulging rule assists in selecting linearizing transformations.

**Transforming Proportions**

① logit  $P \rightarrow \text{logit}(P) = \log \frac{P}{1-P}$

- remove the upper & lower boundaries of the scale, spreading out the tail distribution and making resulting quantities symmetric about 0.

② probit  $P \rightarrow \text{probit}(P) = \Phi^{-1}(P)$

- $\text{logit} \approx (1/\sqrt{P}) \times \text{probit}$

**Linear Algebra Review**

Showing  $H$  is PSD:

(i) Let  $\vec{a} \in \mathbb{R}^n$ . Show  $\vec{a}^T H \vec{a} \geq 0$ .

Norm:  $\|\vec{x}\|_p = (\sum_i x_i^p)^{1/p}$

Dot product  $\langle \vec{x}, \vec{y} \rangle = \vec{x} \cdot \vec{y} = \vec{x}^T \vec{y} = \sum x_i y_i$

$\|\vec{x}\|_2^2 = \vec{x} \cdot \vec{x} = \|\vec{x}\|_2 \|\vec{x}\|_2 \cos 0$

Orthogonality = dot product  $\langle \vec{x}, \vec{y} \rangle = 0$

$\|\vec{x} + \vec{y}\|_2 = \langle \vec{x} + \vec{y} \rangle^T (\vec{x} + \vec{y}) = \vec{x}^T \vec{x} + \vec{x}^T \vec{y} + \vec{y}^T \vec{x} + \vec{y}^T \vec{y}$

Orthogonal basis: If  $\vec{v}_1, \dots, \vec{v}_k$  is orthogonal basis of  $V$ , then  $\vec{u} = \sum c_i \vec{v}_i$  where  $c_i = \vec{u} \cdot \vec{v}_i$ .

Projections:  $\text{argmin}_{\vec{c}} \|\vec{u} - \vec{c}\|_2^2 = \frac{\vec{u}^T \vec{v}_i}{\vec{v}_i^T \vec{v}_i} = c^*$

projecting a vector onto the basis is the minimizer.

ex) project  $\vec{y}$  onto  $\vec{z}$ :  $c^* = \frac{\vec{y}^T \vec{z}}{\vec{z}^T \vec{z}} = \vec{y}$

Gram-Schmidt for generic  $\{\vec{v}_1, \vec{v}_2, \vec{v}_3\}$

(1)  $\vec{u}_1 = \vec{v}_1$

(2)  $\vec{u}_2 = \vec{v}_2 - \text{Proj}_{\vec{u}_1}(\vec{v}_2)$

(3)  $\vec{u}_3 = \vec{v}_3 - \text{Proj}_{\vec{u}_1}(\vec{v}_3) - \text{Proj}_{\vec{u}_2}(\vec{v}_3)$

ex)  $\vec{z} \cdot \vec{y}$  (1)  $\vec{u}_1 = \vec{1}$  (2)  $\vec{u}_2 = \vec{x} - \text{Proj}_{\vec{u}_1}(\vec{x}) = \vec{x} - \vec{1}$

Forsquare Matrix  $A$ , trace( $A$ ) is the sum of diagonal elements.

$\text{tr}(A-B) = \text{tr}(A) - \text{tr}(B)$

**Multiple Regression**

$y_i = \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$

$\vec{y} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \vec{\beta} = \vec{X} \vec{\beta} + \vec{\epsilon}$

Assume  $X$  is full rank  $\Rightarrow \vec{X} \vec{\beta}$  invertible

optimal  $\vec{\beta}$

$$\vec{\beta} = \text{proj}_{\text{span}(\vec{X})}(\vec{y}) + (\vec{y} - \text{proj}_{\text{span}(\vec{X})}(\vec{y})) = \vec{X} \vec{\beta} + \vec{\epsilon}$$

$$X^T \vec{\beta} = K \vec{X} \vec{\beta} + X^T \vec{\epsilon} \Rightarrow 0 \quad (\vec{X} \perp \vec{\epsilon})$$

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

probabilistic viewpoint

consider vector valued random variable  $V = \begin{pmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{pmatrix}$

For scalar  $V$ :  $\text{Var}(V) = \mathbb{E}[V^2] - \mathbb{E}[V]^2$

For vector  $V$ :  $\mathbb{E}[V] = \begin{bmatrix} \mathbb{E}[V_1] \\ \vdots \\ \mathbb{E}[V_n] \end{bmatrix}$   $\text{Var}(V) = \mathbb{E}[(V - \mathbb{E}[V])(V - \mathbb{E}[V])^T] = \begin{bmatrix} \text{Var}[V_1] & \text{cov}(V_1, V_2) & \dots \\ \vdots & \ddots & \vdots \\ \text{cov}(V_n, V_1) & \dots & \text{Var}[V_n] \end{bmatrix}$

(Quadratic Form)  $\vec{V}^T \vec{A} \vec{V} : \vec{V} \in \mathbb{R}^n, \vec{A} \in \mathbb{R}^{n \times n}$  constant, square matrix

- i)  $\vec{V}^T \vec{A} \vec{V} \in \mathbb{R}$
- ii)  $\vec{V}^T \vec{A} \vec{V} = \vec{V}^T \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \vec{V} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} v_i v_j = \sum_{i=1}^n \sum_{j=1}^n a_{ij} v_i v_j$

Some facts about quadratic form

- i) Quadratic Expectation:  $\mathbb{E}[\vec{V}^T \vec{A} \vec{V}] = \text{trace}(\vec{A} \sum_{i=1}^n \vec{V}_i \vec{V}_i^T) + (\mathbb{E}[\vec{V}])^T \vec{A} \mathbb{E}[\vec{V}]$
- ii) If  $\vec{U} = \vec{A}\vec{V}, \vec{W} = \vec{B}\vec{V}$  with  $A, B$  constant, and  $\vec{V}$  random,  $\mathbb{E}[\vec{U}, \vec{W}] = E[(\vec{U} - \mathbb{E}[\vec{U}])(\vec{W} - \mathbb{E}[\vec{W}])^T] = A \Sigma_{i=1}^n \vec{V}_i \vec{V}_i^T B^T$

Model  $\vec{y} \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}, \vec{B} \in \mathbb{R}^p, \vec{\epsilon} \in \mathbb{R}^n$

$$\vec{y} = \vec{X} \vec{B} + \vec{\epsilon}$$

$\mathbb{E}[\vec{y} | X] = \mathbb{E}[\vec{X} \vec{B} + \vec{\epsilon} | X] = \mathbb{E}[\vec{X} \vec{B} | X] + \mathbb{E}[\vec{\epsilon} | X] = X \mathbb{E}[\vec{B} | X] = X \vec{B}$

$\mathbb{E}[\vec{y}, \vec{y}] = \mathbb{E}[(\vec{y} - \mathbb{E}[\vec{y}](\vec{y} - \mathbb{E}[\vec{y}])^T)] = \mathbb{E}[(\vec{X} \vec{B} + \vec{\epsilon} - \mathbb{E}[\vec{X} \vec{B} + \vec{\epsilon}](\vec{X} \vec{B} + \vec{\epsilon})^T)] = \mathbb{E}[\vec{\epsilon} \vec{\epsilon}^T] = \sigma^2 \mathbb{I}_n$

$\mathbb{E}[\vec{B} | X] = \mathbb{E}[(\vec{X} \vec{B})^T \vec{X} \vec{B} | X] = (\vec{X} \vec{B})^T \mathbb{E}[\vec{X} | X] = (\vec{X} \vec{B})^T \vec{X} \vec{B} = B$

$\Sigma_{\vec{B}, \vec{B}} = \mathbb{E}[(\vec{B} - \mathbb{E}[\vec{B}]) (\vec{B} - \mathbb{E}[\vec{B}])^T] = \mathbb{E}[(\vec{X} \vec{B} + \vec{\epsilon} - \mathbb{E}[\vec{X} \vec{B} + \vec{\epsilon}]) (\vec{X} \vec{B} + \vec{\epsilon})^T] = \mathbb{E}[(\vec{X} \vec{B})^T \vec{X} \vec{B} + \vec{\epsilon}^T \vec{\epsilon} - \mathbb{E}[\vec{X} \vec{B}]^T \vec{X} \vec{B} - \mathbb{E}[\vec{\epsilon}]^T \vec{\epsilon}] = \mathbb{E}[\vec{X} \vec{X}^T \vec{B} \vec{B}^T + \vec{\epsilon} \vec{\epsilon}^T] - \mathbb{E}[\vec{X} \vec{B}]^T \vec{X} \vec{B} - \mathbb{E}[\vec{\epsilon}]^T \vec{\epsilon} = \sigma^2 \mathbb{X}^T \mathbb{X} + \sigma^2 \mathbb{I}_n$

$\vec{B}$  under ground truth model:  $N(B, \Sigma_{\vec{B}, \vec{B}})$

In practice, we estimate  $\sigma^2$  as  $\hat{\sigma}^2 = \frac{\sum e_i^2}{n-p-1}$  where  $p$  is # of covariates.

**Ridgeless regression**

**Adjusted R-squared**

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

**Omitted Variable Bias**

Def:  $\vec{y} = \vec{X} \vec{B} + \vec{\epsilon} + \vec{\delta} + \vec{\epsilon}'$ ;  $\vec{X} \in \mathbb{R}^{n \times p}, \vec{\delta} \in \mathbb{R}^n$

(contounder)  $\vec{\delta}$  is a contounder if i) it determines  $\vec{y}$  ii) correlates with other variables of interest  $\vec{X}$ .

OLS without  $\vec{\delta}$ :  $\vec{B} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$

$$= (\vec{X}^T \vec{X})^{-1} \vec{X}^T (\vec{X} \vec{B} + \vec{\epsilon} + \vec{\delta} + \vec{\epsilon}')$$

$$= \vec{B} + (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{\delta} + (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{\epsilon}'$$

**Omitted Variable Bias (OVB)**

If  $X$  and  $\vec{\delta}$  have higher correlation then  $E[X^T \vec{\delta} | X]$  increased

If  $\vec{\delta}$  and  $\vec{X}$  have higher dependence, then  $\vec{\delta}$  increases  $\rightarrow$  bias from omitting  $\vec{\delta}$  is larger.

$MSE := \frac{1}{n} \|\vec{y} - \hat{\vec{y}}\|^2$

partial correlation: linear relationship between  $\vec{y}$  and  $\vec{x}_i$  in the space orthogonal to  $\vec{x}_2$ .

partial correlation between  $\vec{y}$  and  $\vec{x}_i$  given  $\vec{x}_2$  is computed as follows:

- i) regres  $\vec{y}$  on  $\vec{x}_2$ , yielding  $\hat{\vec{y}}_2$ .
- ii) regres each column of  $\vec{x}_1$  on  $\vec{x}_2$ , yielding  $\hat{\vec{x}}_{11}, \dots, \hat{\vec{x}}_{1q}$
- iii) the correlation between  $\vec{y}$  and  $\hat{\vec{x}}_{11}, \dots, \hat{\vec{x}}_{1q}$  are partial correlations.

multiple regression yield parameters that directly correspond to the partial correlations.

**FWL Theorem**

- (i) reducing multivariate regression to univariate ones.
- (ii) Model:  $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$
- $B_1$  are equivalent to:
- (i) OLS estimator obtained by regressing  $y$  on  $x_1$  and  $x_2$
- (ii) OLS estimator obtained by regressing  $y$  on the residual from regression of  $x_1$  on  $x_2$
- (iii) OLS estimator obtained by regressing residual from the regression  $y$  on  $x_2$  on  $x_1$ .

$\vec{y} = \vec{x}_1 \vec{B}_1 + \vec{x}_2 \vec{B}_2 + \vec{\epsilon}$

Let  $M_{X_1} = I_{n \times n} - \vec{x}_1 \vec{x}_1^T$  be the orthogonal complement of  $\vec{x}_1$  projection matrix. Then for OLS regression:

$M_{X_1} \vec{y} = (M_{X_1} \vec{y}) \vec{B}_2 + \vec{\epsilon}$  where  $\vec{B}_2 = \vec{B}$

$\rightarrow \vec{B}_2$  is applied to the part of  $\vec{y}$  that is uncorrelated with  $\vec{x}_1$ .

terms  $M_{X_1} \vec{y}$ : residuals of  $\vec{y}$  on  $\vec{x}_1$ ,  $M_{X_1} \vec{x}_2$ : residuals of  $\vec{x}_2$  on  $\vec{x}_1$ .

$SE^2(\vec{B}_2) = \frac{\sigma^2}{\sum_{i=1}^n (M_{X_1} \vec{x}_i)^2}$

$\rightarrow SE^2(\vec{B}_2) = \frac{\sigma^2}{\sum_{i=1}^n ((M_{X_1} \vec{x}_i) - (M_{X_1} \vec{x}_1))^2} = \frac{\sigma^2}{\sum_{i=1}^n \vec{x}_i^T M_{X_1} \vec{x}_i} = \frac{\sigma^2}{\sum_{i=1}^n \vec{x}_i^T \vec{x}_i - \sum_{i=1}^n \vec{x}_1^T \vec{x}_i}$

$\vec{B}$  is essentially partial correlation with size/variance of  $\vec{x}_1, \vec{x}_2$  into account.

Note that  $\vec{B}$ 's diagonal elements are less than 1. Hence,  $E[\vec{B}_1 \vec{B}_1^T] \neq \sigma^2 \mathbb{I}_2$  as the  $\vec{B}_1$ 's have deflated variance to  $\sigma^2$  ( $= \text{Var}(\vec{B}_1) \leq \sigma^2 \mathbb{I}_2$ ). since there are only  $n-(p+1)$  free components in  $\vec{B}_1$ , this suggests that we can represent  $\vec{B}_1$  as an  $n-(p+1)$  size vector.

**Gram-Schmidt** procedure for error:

Let  $\vec{v} = \text{col}(\vec{x})$

- i) Find basis  $\vec{u}_i$  for  $\vec{v}^\perp = \text{Null}(\vec{x}^T)$  is one rank( $\vec{x}$ ) =  $p+1$ ,  $\vec{G} \in \mathbb{R}^{(n-p-1) \times (n-p-1)}$
- ii) Make  $\vec{G}$  orthogonal using Gram-Schmidt
- iii) Normalize  $\left[ \begin{array}{c} \vec{u}_1 \\ \vdots \\ \vec{u}_{n-p} \end{array} \right]$  s.t.  $\|\vec{u}_i\|^2 = 1$

**Properties of G matrix**

- i)  $\vec{G} \vec{G}^T = \vec{0}$
- ii)  $\vec{G}^T \vec{G} = \vec{I}_{n-p-1}$
- iii) For  $\vec{z}_{p+1} = \vec{G} \vec{v}_1$   $\vec{z}_{p+1} \vec{z}_{p+1}^T = \vec{G} \vec{v}_1 \vec{v}_1^T \vec{G}^T = \vec{G} \vec{G}^T = \vec{0}$

**Table 9.1** Comparison Between Simple Regression Using Scalars and Multiple Regression Using Matrices

	Simple Regression	Multiple Regression
Model	$y = \alpha + \beta x + \epsilon$	$y = X\beta + \epsilon$
Least-squares estimator	$B = \frac{\sum x^2}{\sum x^2} = (\sum x^2)^{-1} \sum x^2 Y^*$	$b = (X'X)^{-1} X'Y$
Sampling variance	$V(B) = \frac{\sigma_e^2}{\sum x^2}$	$V(b) = \sigma_e^2 (X'X)^{-1}$
Distribution	$B \sim N(\beta, \sigma_e^2 (\sum x^2)^{-1})$	$b \sim N_{k+1}[\beta, \sigma_e^2 (X'X)^{-1}]$

NOTE: Subscripts are suppressed in this table; in particular,  $x^* = x_i - \bar{x}$  and  $Y^* = Y_i - \bar{Y}$ .

## Maximum Likelihood Estimation

• Under the linear model,  $y \sim N(X\beta, \sigma_e^2 I_n)$ ,  $\beta$  is MLE of  $\beta$ .

• For the  $i$ th observation,  $y_i \sim N(x_i^* \beta, \sigma_e^2)$  where  $x_i^*$  is the  $i$ th row of model matrix  $X$ .

• PDF of observation  $i$  is:

$$p(y_i) = \frac{1}{\sigma_e \sqrt{2\pi}} \exp \left[ -\frac{(y_i - x_i^* \beta)^2}{2\sigma_e^2} \right]$$

Because the  $n$  observations are independent, their joint probability density is the product of their marginal densities:

$$p(y) = \frac{1}{(\sigma_e \sqrt{2\pi})^n} \exp \left[ -\frac{\sum (y_i - x_i^* \beta)^2}{2\sigma_e^2} \right] = \frac{1}{(2\pi\sigma_e^2)^{n/2}} \exp \left[ -\frac{(y - X\beta)'(y - X\beta)}{2\sigma_e^2} \right] \quad (9.12)$$

• Log likelihood is:  $\log_e L(\beta, \sigma_e^2) = -\frac{n}{2} \log_e 2\pi - \frac{n}{2} \log_e \sigma_e^2 - \frac{1}{2\sigma_e^2} (y - X\beta)'(y - X\beta)$

• Partial derivative is:

$$\frac{\partial \log_e L(\beta, \sigma_e^2)}{\partial \beta} = -\frac{1}{2\sigma_e^2} (2X'X\beta - 2X'y)$$

$$\frac{\partial \log_e L(\beta, \sigma_e^2)}{\partial \sigma_e^2} = -\frac{n}{2} \left( \frac{1}{\sigma_e^2} \right) + \frac{1}{2\sigma_e^4} (y - X\beta)'(y - X\beta)$$

• Der = 0:  $\hat{\beta} = (X'X)^{-1} X'y$ , (MLE  $\hat{\beta}$  = OLS  $\hat{\beta}$ )

$$\hat{\sigma}_e^2 = \frac{(y - X\hat{\beta})'(y - X\hat{\beta})}{n} = \frac{e'e}{n}$$
 (Biased)

• MLE's variance  $\hat{\sigma}_e^2$  is biased. So, we prefer unbiased estimator  $S_E^2 = e'e \rightarrow \hat{\sigma}_e^2$ . MLE is still consistent b/c as  $n \rightarrow \infty$ ,  $\hat{\sigma}_e^2 \rightarrow 0$ .

## Inference for Individual Coefficients

We saw that the least-squares estimator  $\beta$  follows a normal distribution with expectation  $\beta$  and covariance matrix  $\sigma_e^2 (X'X)^{-1}$ .<sup>33</sup> Consequently, an individual coefficient  $\beta_j$  is normally distributed with expectation  $\beta_j$  and sampling variance  $\sigma_{\beta_j}^2$ , where  $v_{jj}$  is the  $j$ th diagonal entry of  $(X'X)^{-1}$ .<sup>34</sup> The ratio  $(\beta_j - \beta_j)/\sigma_{\beta_j}$ , therefore, follows the unit-normal distribution  $N(0, 1)$ , and to test the hypothesis  $H_0: \beta_j = \beta_j^{(0)}$ , we can calculate the test statistic

$$Z_0 = \frac{\beta_j - \beta_j^{(0)}}{\sigma_{\beta_j} \sqrt{v_{jj}}}$$

unknown

$\hat{\sigma}_{\beta_j}^2$  is unknown, but  $S_E^2 = \frac{e'e}{n-p-1}$  can estimate it.

$$\hat{\sigma}_{\beta_j}^2 = S_E^2 (X'X)^{-1} = \frac{e'e}{n-p-1} (X'X)^{-1}$$

• Testing Hypothesis  $H_0: \beta_j = \beta_j^{(0)}$ , we calculate:  $t_0 = \beta_j - \beta_j^{(0)}$

• (WLLN-a): CI for  $\beta_j$ :  $\beta_j = \beta_j \pm t_{1-\alpha/2, n-p-1} S_E(\beta_j)$

The estimated covariance matrix of the least-squares coefficients is  $\hat{V}(\beta) = S_E^2 (X'X)^{-1}$ . The standard errors of the regression coefficients are the square-root diagonal entries of this matrix. Under the assumptions of the model,  $(\beta_j - \beta_j)/S_E(\beta_j) \sim t_{n-p-1}$ , providing a basis for hypothesis tests and confidence intervals for individual coefficients.

## Inference for Multiple Coefficients

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$  where  $L_{H_0}(\bar{y})$  is maximum of the data under the null model,

and the test is rejects  $H_0$  if  $\lambda_{LRT} < c$ .

Thm (Neyman-Pearson Lemma) The LRT is the most powerful test for any given hypothesis. Highest power for a given Type I error rate, i.e. reject when  $H_A$  is true.

$$L_{H_0}(\bar{y}) = \frac{1}{(2\pi\sigma_e^2)^{n/2}} \exp \left( -\frac{(y - X\beta_{MLE})^T (y - X\beta_{MLE})}{2\sigma_e^2} \right) \frac{e'e}{n} = \hat{\sigma}_e^2 \text{ FORME.}$$

$$L_{H_0}(\bar{y}) = C (t + \bar{e}' \bar{e}^T \bar{e}' \bar{e})^{-1/2} = C (\frac{1}{n} RSS_{SS})^{-1/2}$$

$$LRT \text{ is thus, } \lambda_{LRT} = \frac{L_{H_0}(\bar{y})}{L_{H_0}(\bar{y}')} = \frac{(\bar{e}' \bar{e})^{-1/2} RSS}{(\bar{e}' \bar{e}^T \bar{e}' \bar{e})^{-1/2} RSS} = (\lambda_{LRT})^{2/n}$$

\* Fact: If  $\lambda$  is an LRT, and  $f(\lambda)$  is monotone, then  $f(\lambda)$  is an LRT.

Use the fact above to know that F is an LRT for the  $H_0$ .

Recall that for F-test:

$$F_{q, n-p-1} = \frac{(RSS_q - RSS)/q}{RSS/(n-p-1)}$$

$$\frac{q}{n-p-1} F_{q, n-p-1} = \frac{RSS_q - RSS}{RSS} = \frac{RSS_q}{RSS}$$

$$\Rightarrow \left( \frac{q}{n-p-1} F_{q, n-p-1} \right)^q = \frac{RSS_q}{RSS} = \lambda_{LRT}$$

$$\Rightarrow \text{invert} \frac{q}{n-p-1} \left[ \left( \lambda_{LRT, 2} \right)^q - 1 \right] = F_{q, n-p-1}$$

Thus,  $F_{q, n-p-1}$  is an LRT for this hypothesis.

Ho is True	Ha is True		
Reject Ho	Type I Error False Positive $\alpha$ (alpha)	Correct True Positive $1-\beta$ (power)	
Fail to reject Ho	Correct True Negative $1-\alpha$	Type II Error False Negative $\beta$ (beta)	

• BONCH of F-test over LRT

• we know exact dist of  $F_q \{ F(q, n-p-1) \}$

•  $\lambda_{LRT}$  has asymptotic dist.

•  $\sim 2 \ln \lambda_{LRT} \xrightarrow{n \rightarrow \infty} \chi^2_q$

Thus,  $F_{q, n-p-1}$  is an LRT for this hypothesis.

## General Linear Hypothesis

Even more generally, we can test the linear hypothesis

$$H_0: \mathbf{L} \begin{pmatrix} q \times k+1 \\ q \times 1 \end{pmatrix} \begin{pmatrix} \beta \\ \epsilon \end{pmatrix} = \mathbf{c}$$

where  $\mathbf{L}$  and  $\mathbf{c}$  contain prespecified constants, and the hypothesis matrix  $\mathbf{L}$  is of full row rank  $q \leq k+1$ . The resulting F-statistic,

$$F_0 = \frac{(\mathbf{L}\mathbf{b} - \mathbf{c})' [\mathbf{L}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{L}']^{-1} (\mathbf{L}\mathbf{b} - \mathbf{c})}{q S_E^2} \quad (9.17)$$

follows an F-distribution with  $q$  and  $n - k - 1$  degrees of freedom if  $H_0$  is true.

To understand the structure of Equation 9.17, recall that  $\mathbf{b} \sim N_{k+1}[\beta, \sigma_e^2 (X'X)^{-1}]$ . As a consequence,

$$\mathbf{L}\mathbf{b} \sim N_q[\mathbf{L}\beta, \sigma_e^2 \mathbf{L}(X'X)^{-1} \mathbf{L}']$$

Under  $H_0$ ,  $\mathbf{L}\beta = \mathbf{c}$ , and thus

$$(\mathbf{L}\mathbf{b} - \mathbf{c})' [\mathbf{L}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{L}']^{-1} (\mathbf{L}\mathbf{b} - \mathbf{c}) / \sigma_e^2 \sim \chi_q^2$$

$$\text{ex) } \mathbf{L}: \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}; \mathbf{c} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \text{ tests } \beta_1 = \beta_2 = 0$$

$$\mathbf{L} = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}; \mathbf{c} = \begin{bmatrix} 0 \end{bmatrix} \text{ tests } \beta_1 + \beta_2 = 0$$

## Joint Confidence Region

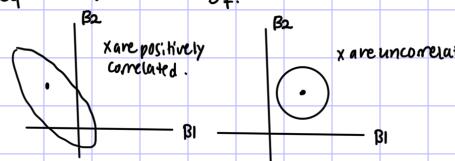
• Inversion of F-test:

1) Let  $\vec{\beta}_q := [\beta_1, \dots, \beta_q]^T$  be our params of interest. If  $\vec{\beta}_q^{(0)}$  is in CI, then we should be able to reject.  $H_0: \vec{\beta}_q = \vec{\beta}_q^{(0)}$  at level  $\alpha$ .

2) If we apply general F-test for  $H_0: \vec{\beta}_q = \vec{\beta}_q^{(0)}$ , we get

$$F_{(1)} = \frac{(\vec{\beta}_q - \vec{\beta}_q^{(0)})' V_q^{-1} (\vec{\beta}_q - \vec{\beta}_q^{(0)})}{q S_E^2(n-p-1)}$$

where  $V_q$  (not inverted) are rows/columns of  $(X'X)^{-1}$  corresponding to  $\vec{\beta}_q$ .



## Dummy Variable Regression

$$x = \begin{pmatrix} \text{small} & \text{large} \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \quad xTx = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{pmatrix} \quad \text{then, } \vec{\beta} = \begin{pmatrix} \beta_0 & 0 \\ 0 & \beta_1 & \beta_2 \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

$$xTy = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ 2\beta_2 \end{pmatrix}$$

An intercept and single dummy gives us the same info.

$$x = \begin{pmatrix} \text{intercept} & \text{large dummy} \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \quad xTx = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{pmatrix} \quad \vec{\beta} = \begin{pmatrix} \beta_0 & \beta_1 & \beta_2 & \beta_3 & \beta_4 \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix}$$

$$xTy = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ 2\beta_4 \end{pmatrix}$$

Interaction terms

$$\bar{y} = \beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 + d(\bar{x}_1 \bar{x}_2) + \epsilon$$

Interaction allows different slope with respect to the continuous covariate.

$$\text{B0} \text{ is slope} = \beta_1$$

$$\text{B1} \text{ is slope} = \beta_1 + d$$

$$\text{Categorical Variables}$$

$$x = \begin{pmatrix} E_1 & E_2 & E_3 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{These columns are orthogonal to each other.}$$

$$\vec{\beta} = \sum_{j=1}^3 \vec{\beta}_j \vec{x}_j = \frac{1}{\| \vec{x}_j \|_2^2} \vec{x}_j \vec{y}_j$$

$$\hat{y}_i = \sum_j \beta_j y_{ij} / n \text{ energy} = \bar{y}$$

with intercept:

$$x = \begin{pmatrix} \text{intercept} & E_2 & E_3 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix} \quad \vec{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

$$\hat{y}_i = \bar{y} \text{ energy}_2 - \bar{y} \text{ energy}_3$$

$$\hat{y}_i = \bar{y} \text{ energy}_3 - \bar{y} \text{ energy}_2$$

If we want to test  $\bar{y}_1 = \bar{y}_2 = \bar{y}_3$ , then we run subset F test against intercept only model.

= principle of marginality

• income by type interaction.

tapply(canada\_occ\$prestige, canada\_occ\$type, mean)

## Analysis of Variance Table

## Model 1: prestige ~ income + education + type + income \* education + type \* education

## Residual Df Sum of Sq F Pr(>F)

## 1 91 4604.6 96.77 1.11 9.92 2.388e-05 \*\*\*

## 2 89 3739.9 2.227 14.511 < 2e-16 \*\*\*

## 3 87 3635.9 2.444 2.748 0.00718 \*\*

## 4 85 3635.9 2.444 2.748 0.00718 \*\*

## Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' '

## Residuals: Min 1Q Median 3Q Max

## -18.273 -7.1773 -0.0854 6.1174 25.2656

##

## Coefficients:

## Estimate Std. Error t value Pr(>|t|)

## (Intercept) 35.527 1.432 24.810 < 2e-16 \*\*\*

## typeprof 32.321 2.227 14.511 < 2e-16 \*\*\*

## typewpc 6.716 2.444 2.748 0.00718 \*\*

## Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' '

## Residual standard error: 9.499 on 95 degrees of freedom

## Multiple R-squared: 0.6976, Adjusted R-squared: 0.6913

## F-statistic: 109.6 on 2 and 95 DF, p-value: < 2.2e-16

→ In categorical variable, if you want to test whether the variable should be in the model or not, then we examine several coefficients. Thus, incremental F-test is useful.

→ BONCH of F-test over LRT

→ we know exact dist of  $F_q \{ F(q, n-p-1) \}$

→  $\lambda_{LRT}$  has asymptotic dist.

→  $\sim 2 \ln \lambda_{LRT} \xrightarrow{n \rightarrow \infty} \chi^2_q$

→  $\lambda_{LRT, 2}$  is an LRT for the  $H_0$ .

→ Use the fact above to know that F is an LRT for the  $H_0$ .

→ Recall that for F-test:

$F_{q, n-p-1} = \frac{(RSS_q - RSS)/(q-1)}{RSS/(n-p-1)}$

→  $\frac{q}{n-p-1} F_{q, n-p-1} = \frac{RSS_q - RSS}{RSS} = \frac{RSS_q}{RSS}$

→  $\left( \frac{q}{n-p-1} F_{q, n-p-1} \right)^q = \frac{RSS_q}{RSS} = \lambda_{LRT}$

→  $\text{invert} \frac{q}{n-p-1} \left[ \left( \lambda_{LRT, 2} \right)^q - 1 \right] = F_{q, n-p-1}$

→ Thus,  $F_{q, n-p-1}$  is an LRT for this hypothesis.

## Bootstrapping

### Nonparametric for $\beta_j$

→ resample our actual data randomly with replacement.

→ compute  $\hat{\beta}_{j,b}$  ( $\hat{\beta}_j$  for sample b)

→ repeat B times, yielding  $\{\hat{\beta}_{j,1}, \dots, \hat{\beta}_{j,B}\}$

→  $\{\hat{\beta}_{j,1}, \dots, \hat{\beta}_{j,B}\}$  is empirical distribution of the sample statistics  $\beta_j$ .

→ Based on the assumption that the sample has approximately the same distribution as the population.

### Parametric Bootstrap [assumes 1) linearity in X, 2) Nomal error]

① run OLS on original data, yielding  $\hat{\beta}$ . If model is correct then  $\hat{y} \sim N(\hat{X}\hat{\beta}, \hat{\sigma}_e^2 I_n)$

② generate n random errors from  $N(0, \hat{\sigma}_e^2 I_n) := \hat{\epsilon}$

$$\hat{y}_b = \hat{y} + \hat{\epsilon}_b$$

③ regress  $\hat{y}_b$  on  $X$  to yield  $\hat{\beta}_{j,b}$

④ Repeat B times to yield  $\{\hat{\beta}_{j,1}, \dots, \hat{\beta}_{j,B}\}$

With these things in mind, let's review the bootstrap hypothesis testing recipe for the case of the F-test.

1. Specify  $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$  or equivalently  $L\beta = \beta_0 = \mathbf{0}_k$ .

2. Choose a test statistic

$$F(\hat{\beta}, \beta_0) = \frac{\frac{1}{k} \{(\hat{\beta} - \beta_0)^T [L(X^T X)^{-1} L]^{-1} (\hat{\beta} - \beta_0)\}}{\frac{1}{n-p-1} ErrSS_{full}}$$

that is large when the observed data deviates far from what is expected under  $H_0$ .

3. Draw bootstrap samples  $\mathbf{y}_1^*, \dots, \mathbf{y}_B^*$  from the sample with replacement. You can do this either by bootstrapping pairs (in which case you will also get bootstrapped design matrices  $X^*$ ) or by bootstrapping residuals.

4. Fit the regression in each bootstrap sample and obtain bootstrap coefficient vectors  $\hat{\beta}_1^*, \dots, \hat{\beta}_B^*$ .

5. Compute bootstrap F-statistics:

$$F(\hat{\beta}_j^*, \hat{\beta}) = \frac{\frac{1}{k} \{(\hat{\beta}_j^* - \hat{\beta})^T [L(X^{*T} X^*)^{-1} L]^{-1} (\hat{\beta}_j^* - \hat{\beta})\}}{\frac{1}{n-p-1} ErrSS_{full,j}} \quad (3)$$

for bootstrap samples  $j = 1, \dots, B$ . Notice that we have followed the bootstrap principle and substituted hat-star quantities for hat quantities (as well as  $X^*$  for  $X$  if bootstrapping pairs) and hat quantities for parameters; this results in a test statistic that looks like it is testing the general linear hypothesis  $L\hat{\beta}_j^* = \hat{\beta}$ , although it's really just the bootstrap version of the usual incremental F-test..

6. Compare the observed value of the test statistic  $F$  to the bootstrap distribution of the test statistic  $F_1^*, \dots, F_B^*$  computed in the previous step. P-values can be computed directly by counting the proportion of bootstrap values that exceed the observed value. Note that although we're using F-statistics for convenience we never look at an F-distribution but rely on the bootstrap to get the sampling distribution.

**Solution:** To test for whether the observation is an outlier, we use the fact that the studentized residuals of a linear regression model are  $t$ -distributed under the standard normal error assumptions. First we construct the studentized residual for observation 42.

$$e_{42}^* = \frac{e_i}{s_{e_{(42)}} \sqrt{1 - h_{42}}}$$

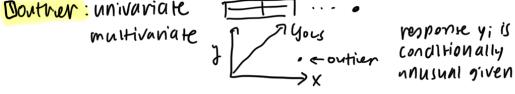
The raw residual  $e_i$  is given by the output from the first regression model; clearly it is the largest negative residual, which is listed as -18.4788. Plugging in the leverage  $h_{42} \approx 0.68$  and the residual standard error from the regression not involving observation 42 — 5.406 from the second regression output — we have

$$e_{42}^* = -18.7488 / (5.406 \sqrt{0.32}) \approx -6.1309$$

This studentized residual will be  $t$ -distributed with  $n - p - 2 = 192$  degrees of freedom under the null hypothesis (the same as the error degrees of freedom in the second regression model). Letting  $t_{192}$  be a random variable with this distribution, we could report the two-sided p-value  $2 \cdot P(t_{192} > | -6.1309 |)$ . However, this would not be quite right since we have selected this residual because of its unusually extreme value. To avoid invalid testing due to this selection step, we use a Bonferroni correction and multiply the p-value by the number of data points in the original regression. Since the model has 7 parameters and 193 error degrees of freedom, there are 200 observations. So the correct p-value for this test is approximately  $400 \cdot P(t_{192} > 6.1309)$ .

## Model Diagnostics

(a) Outliers (b) Leverage (c) Influence



residuals don't have equal variances

$$\Rightarrow V(E_i) = \sigma_e^2 (1-h_i)$$

$\sum_i V(E_i) = \sigma_e^2 \neq \text{values along the diagonal are not equal.}$

Fix: standardized residuals:  $E_i^* = \frac{E_i}{S_{E_i} \sqrt{1-h_i}}$

issue:  $E_i$  and  $S_{E_i}$  are not independent,

$$\tilde{E}_i^* = \frac{E_i}{\sigma_e \sqrt{1-h_i}} = \frac{E_i}{\sigma_e^2 \sqrt{1-h_i}}$$

preventing  $E_i^*$  from following t-distribution.

Fix: Studentized residuals.

Suppose, however, that we refit the model deleting the  $i$ th observation, obtaining an estimate  $S_{E_{(-i)}}$  of  $\sigma_e$  that is based on the remaining  $n - 1$  observations. Then the studentized residual

$$E_i^* = \frac{E_i}{S_{E_{(-i)}} \sqrt{1-h_i}} = \tilde{E}_i^* \sqrt{\frac{n-(p+1)-1}{n-(p+1)-2}} \quad (11.1)$$

has an independent numerator and denominator and follows a t-distribution with  $n - k - 2$  degrees of freedom.

$E_i^*$ : standardised residuals.

procedure:

① refit mean-shift model  $n$  times, once for each observation, producing studentized residuals  $E_1^*, \dots, E_n^*$ .

② usually, our interest focuses on the  $E_{\max}^*$ . However, because we have picked the biggest of  $n$  test statistics, it is not legitimate to find p-value of  $E_{\max}^*$ .

③ To address issue ②, we do Bonferroni adjustment to the p-value for the largest absolute  $E_i$ .

If we have only one residual, we setup:  $P(t_{n-p-2})$

If we have  $n$  residuals, then compute  $P(t_{n-p-2})$  compare to  $\frac{p}{n} = 0.05$   $P(t_{n-p-2})$  ( $E^* > E_{\max}^*$ )

Leverage

idea: is  $x_i^T$  (ith row of  $X$ ) unusual compared to other  $x_j^T$ 's?

leverage ( $x_i^T$ ):  $h_{ii} \leftarrow$  ith diagonal of  $H: X(X^T X)^{-1} X^T$   
properties of  $H$ : ① symmetric ② idempotent ③ rank =  $p+1$

$$\text{④ } \frac{1}{n} \leq h_{ii} \leq 1 \quad \text{⑤ } \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{p+1}{n}$$

\* why is leverage a useful quantity

$\hat{y} = \hat{y}^T \rightarrow \hat{y}_i = \sum_{j=1}^n h_{ij} y_j \times h_{ii}$ : weights for how much  $y_i$  determines  $\hat{y}_i$ .

→ if  $h_{ii}$  is large, then the line near  $y_i$  is almost entirely determined by  $y_i$ .

In simple regression: tells us if  $y_i$  is far from  $\bar{x}$  relative to the full data.

$$h_{ii} = \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \rightarrow n h_{ii} = 1 + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

[Influence] we want to measure how much our estimates ( $\hat{y}$  or  $\hat{\beta}$ ) would change if we remove datapoint  $i$ .

$B_{(-i)}$ : OLS estimate with  $i$ th data point removed.

Q: why not we just do influence;  $= \|\hat{\beta} - \hat{\beta}_{(-i)}\|^2$ ?

A: This is treating  $B$  as in the same scale but each

D may have different units.

[Cook's distance]

Cook's  $D_i$  is the F-statistic for testing the "hypothesis" that  $\beta = \beta_{(-i)}$ :

$$D_i = \frac{(\mathbf{b} - \mathbf{b}_{(-i)})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \mathbf{b}_{(-i)})}{(k+1) S_E^2} = \frac{(\hat{y} - \hat{y}_{(-i)})^T (\hat{y} - \hat{y}_{(-i)})}{(k+1) S_E^2}$$

An alternative interpretation of  $D_i$ , therefore, is that it measures the aggregate influence of observation  $i$  on the fitted values  $\hat{y}$ . This is why Belsley et al. (1980) call their similar statistic "DFFITS". Using Equation 11.8,

$$D_i = \frac{E_i^2}{S_E^2(k+1)} \times \frac{h_{ii}}{(1-h_{ii})^2} = \frac{E_i^2}{k+1} \times \frac{h_{ii}}{1-h_{ii}}$$

which is the formula for Cook's D given in Section 11.4.

\* Cook's distance is considered high if it is greater than 0.5 and extreme if it is greater than 1.

An influence measure proposed by Belsley et al. (1980) closely approximates the squared ratio of volumes of the deleted and full-data confidence regions for the regression coefficients:<sup>18</sup>

$$\text{COVRATIO}_i = \frac{1}{(1-h_{ii}) \left( \frac{n-k-2+p+2}{n-k-1} \right)^{k+1}}$$

Observations that increase the precision of estimation have values of COVRATIO that are larger than 1; those that decrease the precision of estimation have values smaller than 1. Look for values of COVRATIO, therefore, that differ considerably from 1.

\* COVRATIO > 1:  $i$  is not extreme

\* COVRATIO < 1:  $i$  is extreme enough that  $\hat{\beta}_i$  shrinks when it is excluded.

Regress  $\hat{y}$  on  $X$ , get  $\hat{e}_n$ .

2. Regress  $\log(\hat{e}_n^2)$  on  $X$  to yield fitted values, which we call  $\hat{g}_1, \dots, \hat{g}_n$ .

$\hat{h}_i := \exp(\hat{g}_i)$ ;  $\hat{W}^{-1} = \begin{bmatrix} h_1 & & \\ & \ddots & \\ & & h_n \end{bmatrix}$ , and analogously,  $\hat{W} = \begin{bmatrix} 1/h_1 & & \\ & \ddots & \\ & & 1/h_n \end{bmatrix}$ . This is because  $\hat{h}_i$  is the estimator for  $\text{Var}(e_i) = \sigma_e^2/w_i^2$ .

$\hat{\beta}_{FGLS} = (X^T \hat{W} X)^{-1} X^T \hat{W} \hat{y}$ . For inference, we will use the plug-in estimator of the standard error,  $se(\hat{\beta}) = \sigma_e^2 (X^T \hat{W} X)^{-1}$ . Since we have incorporated all of  $\hat{h}_i$  into each  $w_i^2$ , it can also suffice to just set  $\hat{\sigma}_e^2 = 1$ .

↳ Double usage of  $X$  in the above estimation is problematic for the bias of  $\hat{\beta}_{FGLS}$ , but the estimator is still consistent and asymptotically efficient if in fact the errors have different variances as function of  $X$ .

## Assumption Diagnostics

our full model  $\hat{y} \sim N(X\beta, \sigma_e^2 I_n)$

Assumptions are 1) Linearity  $E[\epsilon] = 0$  3) Homoscedasticity 2) Independence of  $\epsilon_i$  4) Normality

### ① Meanzero errors

check studentized residuals plotted against: 1) columns of  $X$  3) order of data/time of collection 2) variables not in the model

$$\hat{\epsilon} = 0 \text{ but the assumption is } E[\epsilon_i | x_i] = 0$$

### ② Nonconstant Variance (Heteroscedasticity)

var( $\epsilon_i | x_i$ )  $= \sigma_e^2 \forall i$

check plot residuals vs fitted values (or vs covariates)

Q: Why not plot against  $\hat{y}$ ?

$\hat{y} = \hat{y} + \epsilon$ , and  $\text{corr}(\hat{y}, \epsilon) = \sqrt{p+1}$  It will be odd to study the correlation between  $\hat{y}$  and  $\epsilon$  that are inherently correlated to each other. But, by construction,  $\hat{y}$  and  $\epsilon$  are orthogonal.

Q: Why prefer plot studentized over  $\epsilon_i$ ?

A: the residuals have unequal variance  $\{\sigma_i^2\}$  even though the random errors might not. In essence, we're removing the variance just due to  $\epsilon$  matrix.

Fixes 1) Transformation 2) Nonparametric Bootstrap

3) WLS (or other standard error procedure)

④ Nonlinearity  $\hat{y} = \hat{y}^T (1/e_i^*)$

Check V-Q plot

$$\hat{y}^T (1/e_i^*) : [0, 1] \rightarrow R$$

$\hat{y}^T (1/e_i^*) = 1$  if  $|e_i^*|$  is smallest.

### Probabilistic Q-Q plots

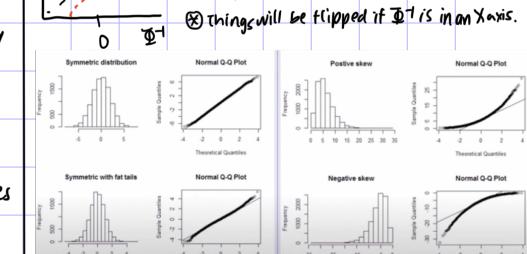
A: right skewed: residuals larger than expected on the right tail

B: left skewed: residuals larger than expected on the left tail

A: large-tailed: more data in the extremes

B: small-tailed: more data in the center

⑧ things will be flipped if  $\hat{y}^T$  is in an x-axis.



Fixes 1) Transformations 2) Bootstrap 3) Assume new error distribution

### Model Selections

Let's say we have  $p$  numeric covariates.

How many parameters for including:

(i) squared covariates:  $p$  more

(ii) pairwise interactions:  $\binom{p}{2}$  more

(iii) three way interactions:  $\binom{p}{3}$  more

Simple/parsimonious is good → interpretable

→ limit future data collection → overfitting concerns.

### Issues with ANOVA/F test

1. Assumption based

2. Significance based tool + predictive importance

3. Paired model comparisons

4. Multiple testing problems

$E[(y_{new} - \hat{y}_{new})^2]$

Model → criterion → quality score

Common criteria: 1) Adj R<sup>2</sup>, 2) cross-validation Error, 3) AIC/BIC, 4) Cp

Each of these criteria help manage the bias/variance tradeoff.

### ① Adjusted R<sup>2</sup>

Ideally, we want to know the error our model will have in predicting unseen data. Adjusted R<sup>2</sup> heuristically penalizes the additional variance that results from extra covariates

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

\* But, adjusted R<sup>2</sup> is not an actual theoretical approximation of  $E[(y_{new} - \hat{y}_{new})^2]$ . It's just a heuristic.

### ② Cross Validation

- cycles through the entire dataset and ensure that every element of  $\hat{y}$  is used as hold out data exactly once.

(i) Leave one out method

$$\text{MSE}(y_{new}, \hat{y}_{new}) = \frac{1}{n} \sum (y_i - \hat{y}_{(-i)})^2 = \frac{1}{n} \sum \frac{(e_i)^2}{1-h_{ii}}$$

(ii) K Folds: hold out random  $K$  fraction of data at a time.

**AIC**

- AIC is based on KL-divergence.
- $D_{KL}(f, f_m) = \int \log\left(\frac{f(y)}{f_m(y)}\right) f(y) dy$
- density under model m.
- Fact:  $D_{KL}(f, f_m) = -E[\log(f_m(y))]$  where c is some function of  $f(y)$
- So if we can estimate  $E[\log(f_m(y))]$ , then we can compare diff models.
- $E[\log(f_m(y))] \approx \frac{1}{n} \sum_i \log(f_m(y_i)) = \frac{1}{n} \log \prod_i f_m(y_i)$
- Approximations. Suggests that  $\min_{f_m} KL \approx \max_{f_m} \text{likelihood}$ .
- We used the same data to approximate the initial expectations & model likelihood. This overstates the expected log-likelihood of  $f_m$  under  $f$ .
- Positive bias of  $\text{Log } L(\hat{B}_{(m)} | y_1, \dots, y_n)$  is  $P_{(m)} + 1$ .
- Thus,  $\underset{m}{\text{argmin}} \{D_{KL}(f, f_m)\} = \underset{m}{\text{argmin}} \{-E[\log(f_m(y))]\}$
- $= \underset{m}{\text{argmin}} \{-E[\log(f_m(y))]\} \approx \underset{m}{\text{argmin}} \{-\log L(\hat{B}_{(m)} | y_1, \dots, y_n) + (P_{(m)} + 1)\}$
- $\rightarrow AIC = -2 \log L(\hat{B}_{(m)} | y_1, \dots, y_n) + 2(P_{(m)} + 1)$
- $\rightarrow BIC = -2 \log L(\hat{B}_{(m)} | y_1, \dots, y_n) + (\log(n))(P_{(m)} + 1)$  parsimonious.
- Lower AIC, BIC means better fit.

**Tradeoffs of using AIC and/or BIC (vs. cross-validation)**

- Only train each model once.
- Use all the data.
- Are based on likelihood model and can be difficult to compare across model types.

"All-subsets" Approach computes  $2^k$  models (if  $k$  total features). One way to lower the total number of models tested is to use Forward / Backward selection

**Forward**

Intercept

Forward & Backward is greedy search. It doesn't guarantee we pick an optimal model. We can lower variance by inducing some bias.

**Shrinkage methods:**

- Preprocessing: is needed b/c
  - we exclude intercept from variable selection
  - the variables are on different scales.
- Standardization:  $X \rightarrow Z = \begin{bmatrix} \vec{x}_1 - \vec{\mu}_1 \\ \vdots \\ \vec{x}_p - \vec{\mu}_p \end{bmatrix}$
- Now regress  $\vec{y}$  on  $Z$ .
- To shrink the parameters, we place constraints on the least squares optimization

**Ridge**

$$\min_B \frac{1}{n} \|\vec{y} - Z\vec{B}\|^2 + \lambda \sum_i B_i^2$$

i) HAS CLOSED FORM SOLUTION

$$f(\vec{B}) = \|\vec{y} - Z\vec{B}\|^2 + \lambda \|\vec{B}\|_2^2$$

$$= (\vec{y} - Z\vec{B})^T (\vec{y} - Z\vec{B}) + \lambda \|\vec{B}\|_2^2$$

$$= \vec{y}^T \vec{y} - 2\vec{y}^T Z\vec{B} + \vec{B}^T Z^T Z\vec{B} + \lambda \|\vec{B}\|_2^2$$

$$= \vec{y}^T \vec{y} - 2\vec{y}^T Z\vec{B} + \vec{B}^T Z^T Z\vec{B} + \lambda \|\vec{B}\|_2^2$$

$$2\vec{B} = 0 - 2Z^T \vec{y} + 2(Z^T Z + \lambda I_p) \vec{B}$$

$$\Rightarrow 2Z^T \vec{y} = 2(Z^T Z + \lambda I_p) \vec{B}$$

$$\Rightarrow \vec{B} = (Z^T Z + \lambda I_p)^{-1} Z^T \vec{y}$$

2) Larger  $\lambda \rightarrow$  smaller  $\vec{B}$   $\Rightarrow$  smaller  $\text{var}(\vec{B})$

$\text{var}(\vec{B}) = (\vec{Z}^T \vec{Z} + \lambda I_p)^{-1} \vec{Z}^T \text{var}(\vec{y}) [\vec{A}]$

$\text{Bias}(\vec{B}) = E[\vec{B}] - \vec{B}$

$= (\vec{Z}^T \vec{Z} + \lambda I_p)^{-1} \vec{Z}^T E[\vec{Z}^T \vec{y}] - \vec{B}$

$= (\vec{Z}^T \vec{Z} + \lambda I_p)^{-1} \vec{Z}^T \vec{Z} \vec{B} - \vec{B}$

diverges from  $\vec{I}^{-1} p$  as  $\lambda$  gets larger.

\* there is ALWAYS some value  $\lambda > 0$  s.t. the MSE of Ridge regression is lower than OLS.  $\lambda$  is unknown/binary

**PRO**

- Analytical solution for each.
- Theoretically can always output from OLS
- No explicit variable selection.
- No clear inference approach.

**CON**

- No clear inference approach.
- OLS is preferable for causal inference since ridge coefficients have bias and don't come with CI.

**Lasso**

$$\min_B \frac{1}{n} \|\vec{y} - Z\vec{B}\|^2 + \lambda \|\vec{B}\|_1$$

i) LASSO DOESN'T HAVE CLOSED FORM

→ but, since it is still convex, there is a single local optimum. There are fast numerical solvers (coordinate descent, iterative LSS, proximal gradient)

**FIGURE 6.6.** The standardized lasso coefficients on the Credit data set are shown as a function of  $\lambda$  and  $\|\vec{B}\|_1/\|\vec{B}\|_2$ .

**FIGURE 6.7.** Contours of the error and constraint functions for the lasso and ridge regression problems. The red contours represent the error function,  $\|\vec{y} - Z\vec{B}\|_2^2$ , and the blue contours represent the constraint function,  $\|\vec{B}\|_1$  or  $\|\vec{B}\|_2$ , while the red ellipses are the contours of the RSS.

**PRO**

- Variable selection
- Convex/computationally feasible
- Nonlinear solution to analyze
- Still no clear non-Bayesian inference

\* use CV to find the optimal  $\lambda$ .

# of regressors (excluding intercept) with nonzero coefficients in the fitted model associated with corresponding  $\lambda$ .

**Generalized Linear Models (GLM)**

- Framework that allows  $\vec{y}$  to follow various parametric distributions: Bernoulli (y is binary), Poisson, NB, Exponential, Gamma, Multinomial (y: count of occurrences of k different outcomes).
- no need of  $\ln K(f(x))$ .
- GLM recipe

- Specify the conditional distribution of  $y_i | \vec{x}_i^T$  (e.g.  $y_i \sim \text{Bernoulli}(\pi_i)$ )
- Consider parameter of distribution in (1), their relation to  $E[y_i | \vec{x}_i^T]$ , and their domain (e.g.  $y_i | \vec{x}_i^T \sim \text{Ber}(\pi_i)$ ,  $\pi_i = E[y_i | \vec{x}_i^T]; \pi_i \in (0, 1)$ )
- Find link function  $g(\cdot)$  that transforms  $E[y_i | \vec{x}_i^T]$  to the range:  $(-\infty, \infty)$ . (e.g.  $y_i | \vec{x}_i^T \sim \text{Ber}(\pi_i)$ , let  $g(\cdot)$  be logit function  $g(\pi) = \log\left(\frac{\pi}{1-\pi}\right)$ ). Then,  $g(\pi_i) \rightarrow (-\infty, \infty)$ , as desired
- Assume that  $g(\pi_i)$  can be modeled linearly.  $g(\pi_i) = \vec{x}_i^T \vec{\beta}$ . we then define inverse link function  $g^{-1}(\cdot)$ . (e.g.  $g(\pi_i) = \log\left(\frac{\pi_i}{1-\pi_i}\right) = \vec{x}_i^T \vec{\beta}$ )
- $\Rightarrow \pi_i = \exp\left(\vec{x}_i^T \vec{\beta}\right) \quad \pi_i = \exp(\vec{x}_i^T \vec{\beta}) - \exp(\vec{x}_i^T \vec{\beta})$
- $\Rightarrow \pi_i = \frac{\exp(\vec{x}_i^T \vec{\beta})}{1 + \exp(\vec{x}_i^T \vec{\beta})} = \frac{1}{1 + \exp(-\vec{x}_i^T \vec{\beta})}$
- $\therefore g^{-1}(\vec{x}_i^T \vec{\beta}) = g^{-1}(\pi_i)$

**Table 15.2** Canonical Link, Response Range, and Conditional Variance Function for Exponential Families

Family	Canonical Link	Range of $Y_i$	$V(Y_i   \eta_i)$
Gaussian	Identity	$(-\infty, +\infty)$	$\phi$
Binomial	Logit	$0, 1, \dots, n_i$	$\mu_i(1 - \mu_i)$
Poisson	Log	$0, 1, 2, \dots$	$\mu_i$
Gamma	Inverse	$(0, \infty)$	$\phi \mu_i^2$
Inverse-Gaussian	Inverse-square	$(0, \infty)$	$\phi \mu_i^3$

NOTE:  $\phi$  is the dispersion parameter,  $\eta_i$  is the linear predictor, and  $\mu_i$  is the expectation of  $Y_i$  (the response). In the binomial family,  $n_i$  is the number of trials.

Compare OLS with GLM:

- OLS:  $E[y_i] = \vec{x}_i^T \vec{\beta}$
- GLM:  $g(E[y_i]) = \vec{x}_i^T \vec{\beta}$

Now, the variance of  $y_i$  is often tied to its mean so it can fluctuate. ex)  $\text{Ber}(\pi_i)$ 's variance:  $\pi_i(1 - \pi_i)$  POIS( $\lambda_i$ )'s var:  $\lambda_i$

**Table 15.1** Some Common Link Functions and Their Inverses

Link	$\eta_i = g(\mu_i)$	$\mu_i = g^{-1}(\eta_i)$
Identity	$\mu_i$	$\eta_i$
Log	$\log \mu_i$	$e^{\eta_i}$
Inverse	$\mu_i^{-1}$	$\eta_i^{-1}$
Inverse-square	$\mu_i^{-2}$	$\eta_i^{-1/2}$
Square-root	$\sqrt{\mu_i}$	$\eta_i^{1/2}$
Logit	$\log \frac{\mu_i}{1 - \mu_i}$	$\frac{1}{1 + e^{-\eta_i}}$
Probit	$\Phi^{-1}(\mu_i)$	$\Phi(\eta_i)$
Log-log	$-\log_e[-\log_e(\mu_i)]$	$\exp[-\exp(-\eta_i)]$
Complementary log-log	$\log_e[-\log_e(1 - \mu_i)]$	$1 - \exp[-\exp(\eta_i)]$

NOTE:  $\mu_i$  is the expected value of the response,  $\eta_i$  is the linear predictor, and  $\Phi(\cdot)$  is the cumulative distribution function of the standard-normal distribution.

A generalized linear model (or GLM) consists of three components:

- A random component, specifying the conditional distribution of the response variable,  $Y_i$  (for the  $i$ th of  $n$  independently sampled observations), given the values of the explanatory variables in the model. In the initial formulation of GLMs, the distribution of  $Y_i$  was a member of an exponential family, such as the Gaussian, binomial, Poisson, gamma, or inverse-Gaussian families of distributions.
- A linear predictor—that is, a linear function of regressors,
- A smooth and invertible linearizing link function  $g(\cdot)$ , which transforms the expectation of the response variable,  $\mu_i = E(Y_i)$ , to the linear predictor:

$g(\mu_i) = \eta_i = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$

Summary of advantages of GLMs over traditional (OLS) regression

- We do not need to transform the response to have a normal distribution.
- The choice of link is separate from the choice of random component, giving us more flexibility in modeling.
- The models are fitted via maximum likelihood estimation, so likelihood functions and parameter estimates benefit from asymptotic normal and chi-square distributions.
- All the inference tools and model checking that we will discuss for logistic and Poisson regression models apply for other GLMs too: e.g., Wald and Likelihood ratio tests, deviance, residuals, confidence intervals, and overdispersion.
- There is often one procedure in a software package to capture all the models listed above, e.g. PROC GENMOD in SAS or glmer in R, etc., with options to vary the three components.
- The data  $Y_1, Y_2, \dots, Y_n$  are independently distributed, i.e., cases are independent.
- The dependent variable  $Y_i$  does NOT need to be normally distributed, but it typically assumes a distribution from an exponential family (e.g. binomial, Poisson, multinomial, normal, etc.).
- A GLM does NOT assume a linear relationship between the response variable and the explanatory variables, but it does assume a linear relationship between the transformed expected response in terms of the link function and the explanatory variables; e.g., for binary logistic regression  $\text{logit}(x) = \beta_0 + \beta_1 x$ .
- Explanatory variables can be nonlinear transformations of some original variables.
- The homogeneity of variance does NOT need to be satisfied. In fact, it is not even possible in many cases given the model structure.
- Errors need to be independent but NOT normally distributed.
- Parameter estimation uses maximum likelihood estimation (MLE) rather than ordinary least squares (OLS).

**Estimation + the data model for Bernoulli's**

Now  $y_i \sim \text{er}(g^{-1}(\vec{x}_i^T \vec{\beta}))$ . This is no longer a least squares problem (so it doesn't have closed form solutions.) But we can still write the likelihood and  $\log L(\vec{\beta}, \vec{y}; \vec{x})$  is convex for many GLMs and we can solve it if dispersion ( $\phi$ ) must be estimated, we can do alternating MLE.

- Init:  $t = 0, \vec{\beta}_0$
- while not converged:
- $\vec{\beta}_{t+1} = \hat{\vec{\beta}}_{MLE} | \vec{\beta}_t$
- $\vec{\beta}_{t+1} = \hat{\vec{\beta}} | \vec{\beta}_{t+1}$
- $t = t + 1$

**Prediction** Once we optimize the GLM for  $\vec{\beta}$ , prediction is straightforward.

$\vec{y}_i = E[y_i | \vec{x}_i^T \vec{\beta}] = g^{-1}(\vec{x}_i^T \vec{\beta})$  [e.g. logistic regression]

$\vec{y}_i = \frac{1}{1 + \exp(-\vec{x}_i^T \vec{\beta})}$

**Model performance:**

- residual Deviance ( $D_m$ ): model with one parameter for each observation.
- $D_m = 2(\log L_S - \log L_m)$
- $L_S$ : likelihood of saturated model
- $L_m$ : likelihood of our model.

\* Fact: For logistic regression,  $D_m = -2 \log L_m$

upf: In saturated model, we can fit the training log perfectly. (i.e. if  $y_i = 1$ , then  $g^{-1}(\vec{x}_i^T \vec{\beta}_{MLE}) = 1$ )

$\vec{L} = \sum_i \pi_i y_i (1 - \pi_i)^{-1} y_i \in \text{Benoulli likelihood} = 1 \Rightarrow \log \text{likelihood} = 0$

$\rightarrow D_m \approx 0 \Rightarrow D_m = -2 \log L_m$

$\rightarrow D_m \sim \chi^2(n-p+1)$  params of null model.

$\rightarrow D_m$  has  $n-p-1$  df

**Comparing  $D_m$  to Null Deviance ( $D_0$ ):**

- $D_0$  is the residual deviance for an intercept only model
- $\rightarrow D_0 \sim \chi^2(n-p)$

**Natural Exponential Family** [gaussian, Binom, poisson, gamma]

& PDF can be written:

$p(y_i; \theta; \phi) = \exp\left[\frac{y_i \theta - b(\theta)}{a(\phi)} + c(y_i, \phi)\right]$

$\theta := g_c(\mathbb{E}[y])$  is some function of  $\mathbb{E}[y]$  called canonical link.  $b(\theta)$   $c(y_i, \phi)$

**Binomial**

$p(y_i; \theta; \phi) = \exp\left[\frac{y_i \theta - \log(1 + \exp(\theta))}{\phi} + \log\left(\frac{n}{y_i}\right)\right]$

Where  $\theta = \log\left(\frac{p}{1-p}\right)$ ,  $n$  is known,  $\phi$  is irrelevant.

**canonical link is  $\log\left(\frac{p}{1-p}\right)$ .**

- $\mathbb{E}[y_i] = b'(\theta)$ ,  $\text{var}(y_i) = a(\phi)b''(\theta)$
- log likelihood of data point  $y_i$  is:  $\ell_i = y_i \theta_i - b(\theta_i) - c(y_i, \phi)$
- $\ell_i = \frac{y_i \theta_i - b(\theta_i)}{\phi} - c(y_i, \phi)$
- $\ell_i := \vec{x}_i^T \vec{\beta} \Rightarrow g(\vec{\beta}) = \ell_i$
- $\ell_i := g_c(\mathbb{E}[y]) := g_c(g^{-1}(\vec{x}_i^T \vec{\beta}))$

$\rightarrow$   $\mathbb{E}[y_i] = M_i$

$\rightarrow g$  is selected linkfn.

to find  $\vec{\beta}_{MLE}$ , we need to find the critical pt of the log likelihood:

$\frac{\partial \ell_i}{\partial \beta_j} = \frac{x_{ij}}{a(\phi)} \frac{y_i - M_i}{b''(\theta_i)}$  for generically selected link function (.)

Full model likelihood is:

$\log L(\vec{\beta}, \vec{y}; \vec{x}) = \sum_i \ell_i$

thus, log likelihood gradient (score):

$S(\vec{\beta}) = \sum_{i=1}^n \frac{\partial \ell_i}{\partial \beta_0}, \dots, \sum_{i=1}^n \frac{\partial \ell_i}{\partial \beta_p}$

where  $D = \text{diag}\left(\left\{g'(M_i) a_i (\phi) b''(\theta_i)\right\}_{i=1, \dots, n}\right)$