

Lecture 1b 2023 (10.11.25 Fall).

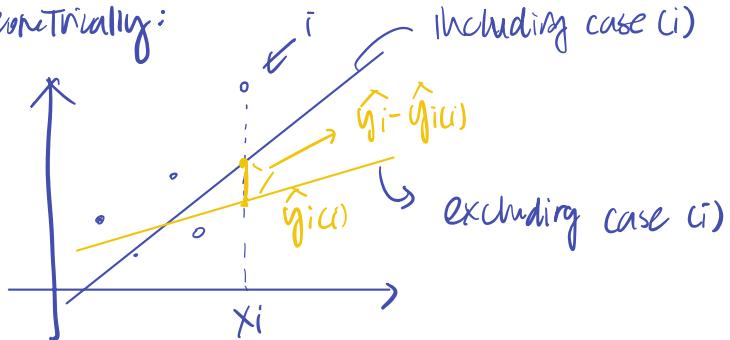
- ① LOOCV DFFITS, DFBETAS, Cook's distance.
- ② Multicollinearity

LOOCV = leave-one-out cross validation

idea: to assess how including/excluding case i affects \hat{y}_i , or b_k , \hat{y}

$$\rightarrow (\text{DFFITS})_i = \frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{\text{MSE}_{(i)} h_{ii}}}$$

Geometrically:



$$\text{Var}(\hat{y}) = \text{Var}(H \cdot y) = H \underbrace{\text{Var}(y)}_{\sigma^2 I_n} H^T = \sigma^2 H H^T = \sigma^2 H$$

$$\text{Var}(\hat{y}_i) = \sigma^2 h_{ii}$$

$$\text{Var}(\hat{y}_i) = \underbrace{\text{MSE}_{(i)}}_{\text{estimated}} h_{ii} \quad \text{exclude case } i \quad \cancel{(\text{MSE}_{(i)})}$$

$\text{MSE}_{(i)}$ might be less bias than MSE .

$$\text{Fact: DFFITS}_i = e_i \left(\frac{(n-p-1)}{\text{SSE}(1-h_{ii}) - e_i^2} \right)^{\frac{1}{2}} \left(\frac{h_{ii}}{1-h_{ii}} \right)^{\frac{1}{2}}$$

$\hookrightarrow t_i$ (extremely studentized residual, studentized e_i) .

Robust linear Regression might be a useful model to use in real world if there are many outliers in the dataset. (not covered in the class).

Rule of thumb: influential on \hat{y}_i

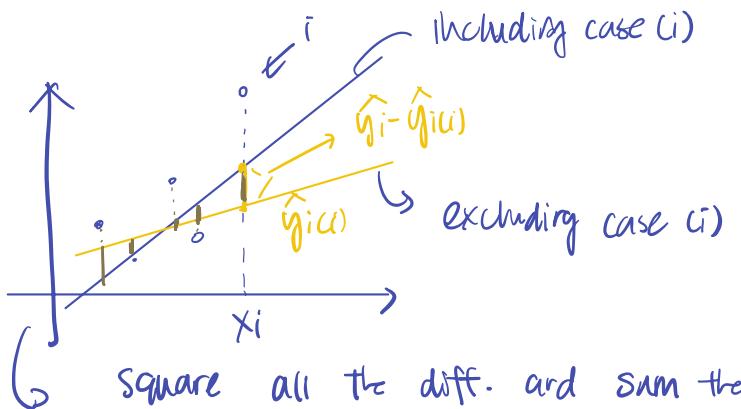
If $|DFFITS_i| > 1$ for "small" / "medium" datasets. Should be # of b .
 $|DFFITS_i| > 2 \frac{P}{n}$ for "large" datasets. $P = \# \text{ of the predictors.}$

$DFFITS_i$ = influence on i th fitted value by the i th case X_i .

Cook's distance := influence on All fitted values by the i th case X_i .

$$D_i = \sum_{j=1}^p (\hat{y}_{ij} - \hat{y}_{j(i)})^2 / (p \cdot \text{MSE}) = \|\hat{y} - \hat{y}_{(i)}\|^2 / \text{RMSE}$$

Cook's Distance.



$$\text{Exercise: } \| \hat{\beta} - \hat{\beta}_{(i)} \|_2^2 = (\underline{b} - \underline{b}_{(i)})^T (X^T X) (\underline{b} - \underline{b}_{(i)}) \quad \leftarrow \text{Check that this is } X_p^2. \quad (*)$$

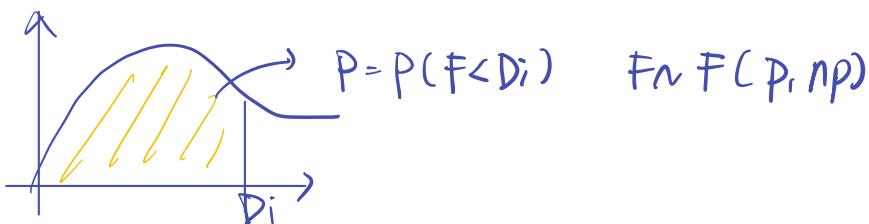
$$\| \hat{\beta} - \hat{\beta}_{(i)} \|_2^2 / \sigma^2 \sim X_p^2 \quad \text{SSE} / \sigma^2 \sim X_{n-p}^2 \quad \text{so: } (X_p^2 / p) / (X_{n-p}^2 / n-p)$$

$\Rightarrow D_i \sim F(p, n-p)$ (*) and that numerator / denominator are $\perp\!\!\!\perp$. (Independent)

Square all the diff. and sum them up.
and then normalized by PMSE to obtain D_i .

Standardize : make the data with mean zero, variance 1.

Normalize : make the data in between 0 & 1.

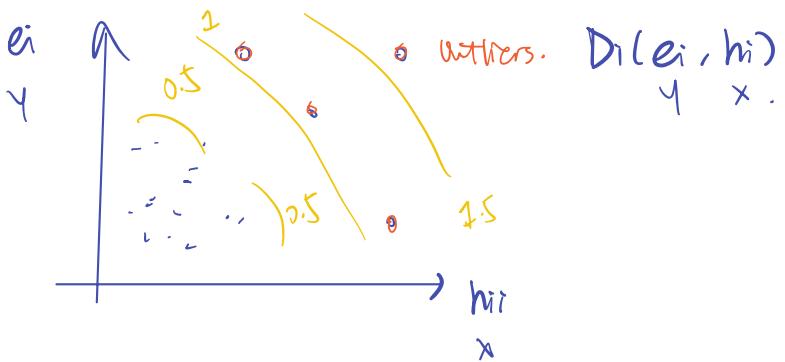


Rule of thumb: $P \in [0.1, 0.2] \leftarrow \text{no influence}$
 $P \in [0.5, 1] \leftarrow \text{major influence.}$

$$\text{Fact: } D_i = (e_i^2 / \text{PMSE}) \cdot (h_{ii} / (1-h_{ii}))^2 = D \underbrace{(e_i, h_{ii})}_{\text{Residual}} \text{ leverage.}$$

Influential case: if 1) e_i is large, h_{ii} is moderate $0 < h_{ii} \leq 1$
 2) e_i is med., h_{ii} large
 3) both are large.

In R: `plot(lm(...))`.



→ Influence on a regression coeff. (of a single case i) .

$$\text{DFBETAS}_{k(i)} = \frac{(b_k - b_{k(i)})}{\sqrt{\text{MSE}_{(i)} C_{kk}}} \quad \text{where } C_{kk} = I(X^T X)^{-1} J_{k+1, k+1}, \quad 0 \dots p-1$$

difference.

= a change in the k th estimated b_k when the i th obsn. x_i is removed, studentized.

$$\text{Var}(b) = \sigma^2 (X^T X)^{-1} \Rightarrow \text{Var}(b_k) = \sigma^2 (I(X^T X)^{-1})_{k+1, k+1}$$

Rule of thumb : $| \text{DFBETAS}_{k(i)} | > 1$, "Small" or "medium" dataset
 $| \text{DFBETAS}_{k(i)} | > 2/\sqrt{n}$, "large" dataset

x_i is influential on b_k

$$\text{NATZ: (Cook's distance)} i \approx \sum_{k=0}^{p-1} (\text{DFBETAS})_{k(i)}^2$$

► Multicollinearity (M.C.)

$$\text{corr}(x_i, \sum_{j \neq i} a_j x_j) \Rightarrow \det(X^T X) \approx 0.$$

Cols X are almost linearly dependent

Causes multicollinearity.

→ Informal diagnostics of M.C.

- ① b_k change a lot when include / exclude predictor x_k .
- ② b_k is not significant for important predictors
- ③ $b_k > 0$ but should be < 0 (or vice versa)
- ④ large Σ_{xx}
- ⑤ wide CI's for b_k

→ Variance inflation factor (VIF) $\propto \text{var of } b_k$'s

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

$$\text{Var}(b^*) = \sigma^2 (\Sigma_{xx})^{-1} \quad (\text{standardized model}).$$

Def: $\text{VIF}_k = [\Sigma_{xx}^{-1}]_{kk}$.

Fact: $VIF_k = \frac{1}{1-R_k^2}$, $k=1, 2, \dots, p-1$.

R_k^2 = coeff. of multiple determination when X_k is regressed on $X_1, \dots, X_{k-1}, X_{k+1}, \dots, X_{p-1}$

$p-2$ predictors (i.e. excluding X_k).

If $R_k^2 = 0$ i.e. X_k is not linearly related with other pred $\Rightarrow VIF_k = 1$.

If $R_k^2 \rightarrow 1 \Rightarrow VIF_k \rightarrow \infty \Rightarrow \text{Var}(b_k^*) \rightarrow \infty$.

Rule of thumb: $VIF_k > 10 \Rightarrow$ bad, M.C. is. present

$$\uparrow \\ R_k^2 \approx 0.9.$$

ISR Internal studentized residual

$$\frac{e_i}{\sqrt{\text{MSE}(1-h_{ii})}}$$

Material -

ESR external stu. res.

$$\frac{e_i}{\sqrt{\text{MSE}_{(i)}(1-h_{ii})}}$$

excluded, external.

* Lecture 18 2023 年 11 月 6 日.

** Weighted LS

$$y = X\beta + \varepsilon, \varepsilon \sim N(0, \sigma^2 I) \quad \mathbb{E}(\varepsilon_i) = 0 \\ \text{Var}(\varepsilon_i) = \sigma^2 \quad \leftarrow \text{heteroskedasticity} \quad \text{cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j \\ \text{Var}(\varepsilon_i) = \sigma_i^2$$

Before: How to fix nonconstant variance? Transform y via Box-Cox transform.
But you may destroy line relation btw. X & y . Var. stabilizing.



messy!
transforms fail to stabilize variance.

Answer: WLS (Weighted Least Square)
as opposed to OLS. (Ordinary Least Square).

Case 1: Assume σ_i^2 's are given

$$L(\beta) = \prod_{i=1}^n \left(\frac{1}{2\pi\sigma_i^2} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2} \sum_{i=1}^n \frac{(y_i - \beta_0 - \beta_1 x_{i,1} - \dots - \beta_p x_{i,p})^2}{\sigma_i^2} \right)$$

↑ likelihood

$$w_i := \frac{1}{\sigma_i^2}, \quad L(\beta) = \prod_{i=1}^n \left(\frac{w_i}{2\pi} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2} \sum_{i=1}^n w_i (y_i - \beta_0 - \beta_1 x_{i,1} - \dots - \beta_p x_{i,p})^2 \right) = Q_W(\beta)$$

Maximizing $L(\beta)$ wrt $\beta \Leftrightarrow$ minimizing $Q_W(\beta)$ wrt β

$$Q_W(\beta) = \sum_{i=1}^n w_i (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_{p-1} x_{ip-1})^2 \\ = (\underline{y} - \underline{x}\beta)^T W (\underline{y} - \underline{x}\beta)$$

where $W = \begin{bmatrix} w_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & w_n \end{bmatrix}$

Observations with high σ_i^2 have low w_i
Observations with low σ_i^2 have high w_i

$\min_{\beta} Q_W(\beta) \rightarrow$ normal equations for WLS.

$$(\underline{x}^T W \underline{x}) \underline{b}_W = \underline{x}^T W \underline{y} \Rightarrow \boxed{\underline{b}_W = (\underline{x}^T W \underline{x})^{-1} \underline{x}^T W \underline{y}}$$

estimated β using WLS.

If $\text{Var}(h_i) = \theta^2 \Rightarrow$ you can show that $\underline{b}_W = \underline{b}^{\text{OLS}}$

Properties : \underline{b}_W = unbiased
= min. variance
= consistent (i.e. $n \rightarrow \infty \Rightarrow \underline{b}_W \rightarrow \beta$)

Case 2 : Variances are known up to a proportionality constant, i.e. we know $\frac{\sigma_i^2}{\theta^2}$
Define $w_i = \frac{k}{\theta_i^2} \rightarrow$ unknown
 $\theta_i^2 \rightarrow$ this is all we know.

Example: happens when response variable is an average of groups of sizes h_i :

$$\underline{y}_{ij} = \underbrace{\frac{1}{h_i} \sum_{k=1}^{h_i} y_{ik}}_{\text{response variable}} \quad \text{constant variance } \theta^2$$

$$\text{Var}(\underline{y}_{ij}) = \frac{\theta^2}{h_i} \rightarrow \text{estimate}$$

$$\Rightarrow \text{Var}(h_i) = \theta^2, \quad \theta^2/\theta^2 = h_i/h_i$$

In this case, $\text{Var}(\underline{b}_W) = k (\underline{x}^T W \underline{x})^{-1}$

$$\hat{s}^2(\underline{b}_W) = \text{MSE}_W (\underline{x}^T W \underline{x})^{-1}$$

$$\text{MSE}_W = \text{estimator of } k = \frac{1}{n-p} \sum_{i=1}^n w_i e_i^2 = \hat{k}_{\text{MLE}}$$

Before: $\text{Var}(\underline{b}) = \theta^2 (\underline{x}^T \underline{x})^{-1}, \quad \hat{s}^2(\underline{b}) = \text{MSE} (\underline{x}^T \underline{x})^{-1}$

Case 3: Variances σ_i^2 are unknown

$$\text{Var}(h_i) = \mathbb{E}(h_i^2) - (\mathbb{E}(h_i))^2 = \mathbb{E}(e_i^2) = \theta_i^2 \approx e_i^2$$

Plot $|e_i|$ or e_i^2 vs. x_{ik} or vs. \underline{y}_{ij} .

$$\frac{\underline{e}^T W \underline{e}}{n-p} \}$$

If $|e_i|$ or e_i^2 plot has shape like , then regress $|e_i|$ or e_i^2 onto X_k wrt X_{ik}

or $\underbrace{X_1 + \dots + X_{p-1}}_{\text{wrt } q_i}$ wrt q_i

- $|e_i| \sim X_{ik} \Rightarrow$ fitted values $|e_i| := \hat{s}_i$ estimate of $\text{sd } \sigma_i$
- $e_i^2 \sim X_1 + \dots + X_{p-1} \Rightarrow$ fitted values of $e_i^2 := \hat{v}_i$ estimate of variance σ_i^2

choose $w_i = \frac{1}{\hat{s}_i^2}$ or $\frac{1}{\hat{v}_i} \Rightarrow b_{\text{WLS}} = (X^T W X)^{-1} X^T W Y$ the WLS estimate
sanity check $MSE_W \approx 1$ since $MSE_W = \hat{R}$

Procedure: IWLs / IRLS
 \downarrow iteratively reweighted LS
 iteratively weighted LS

- Step 1: Fit regression via DLS $\Rightarrow e_i$
- Step 2: Regress $|e_i|$ or e_i^2 onto X_{ik} or q_i
- Step 3: get $w_i = 1/\hat{s}_i^2$ or $w_i = 1/\hat{v}_i$
- Step 4: fit regression using WLS
 If b_{WLS} is significantly different than b
 \hookrightarrow means residual estimates were not good enough.
 \Rightarrow do steps (2)–(4) again
 (now e_i 's in (2) are from (4)) $e_i = y_i - X b_{\text{WLS}}$

Usually 2–3 iterations is enough.

What if we use OLS in case $\sigma_i^2 \neq \text{constant}$?

$$\underline{b} = (X^T X)^{-1} X^T Y \rightarrow \text{unbiased, consistent, not min. Variance}$$

$$\text{Var}(\underline{b}) \neq \sigma^2 (X^T X)^{-1} \quad \hookrightarrow \neq$$

$$= (X^T X)^{-1} (X^T \text{Var}(\underline{e}) X) (X^T X)^{-1}$$

$$\underbrace{\text{Var}(\underline{e})}_{\text{Unknown}} = \begin{bmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_n^2 \end{bmatrix}$$

$$\underline{S}^2(\underline{b}) = (X^T X)^{-1} (X^T \underline{S} \underline{e} X) (X^T X)^{-1} \quad \text{where } \underline{S} \underline{e} = \begin{bmatrix} e_1^2 & \cdots & 0 \\ 0 & \ddots & e_n^2 \end{bmatrix}$$

White's estimator / robust covariance matrix of \underline{b} .

Inference in WLS case:

$$\boxed{S^2(b_{\text{WLS}}) = \text{MSE}_W (X^T W X)^{-1}} \rightarrow \text{CI's, PI's etc.}$$

$$\text{Var}(b_{\text{WLS}}) = 1 \cdot (X^T W X)^{-1}$$

$$S^2(b) = \text{MSE} (X^T X)^{-1}$$

g g sigma

g (97) \$ sigma

Summary

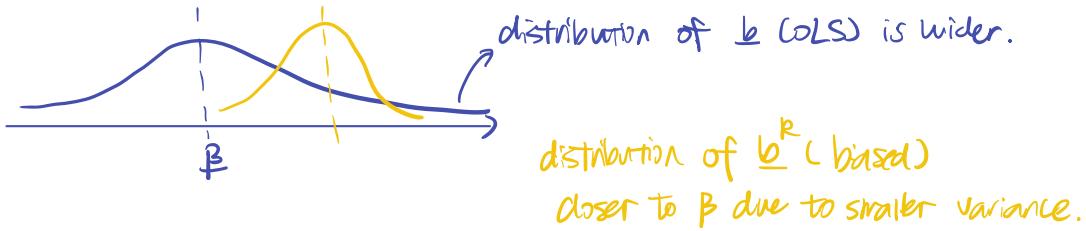
II. Ridge Regression (11.2)

Remedial measure for Multicollinearity.

Signs of M.C.: $\text{sd}(\hat{\beta}_{\text{OLS}})$ are very large \Rightarrow OR
 β_i 's for important pred's not significant.

Idea: reduce variance, sacrifice unbiasedness.

i.e.: finds $\hat{\beta}^R$ s.t. "Var($\hat{\beta}^R$) is small", but $E(\hat{\beta}^R) \neq \beta$



$X_{ik}, Y_i \rightarrow$ Correlation transform

$$\Rightarrow \hat{\beta} = \mathbf{r}_{xx}^{-1} \mathbf{r}_{yx}$$

$$\Rightarrow \text{Ridge regression: } \hat{\beta}^R = (\mathbf{r}_{xx} + c \cdot \mathbf{I})^{-1} \cdot \mathbf{r}_{yx}$$

improving conditioning of the matrix to -1
appropriately chosen constant.

Lecture 19. 2023 11月8日

(I) Ridge Regression

↳ Remedial measure for M.C.

$y_i, X_{ik} \rightarrow$ corr'n transform (entries btw -1, 1)

y_i^*, X_{ik}^*

$$\text{OLS: } \hat{\beta} = \mathbf{r}_{xx}^{-1} \mathbf{r}_{yx}$$

$$\text{Ridge regression: } \hat{\beta}^R = (\mathbf{r}_{xx} + c \cdot \mathbf{I})^{-1} \cdot \mathbf{r}_{yx}$$

appropriately chosen constant

★ Point of view (I): Minimization with regularization.

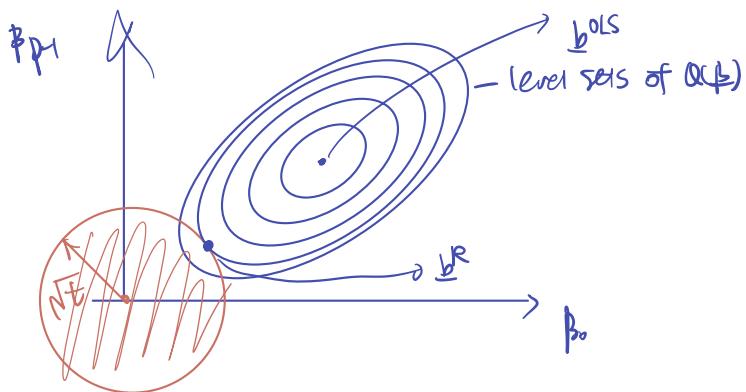
$$\hat{\beta}^R = \underset{\beta}{\operatorname{argmin}} \underbrace{\|y - X\beta\|^2}_{\text{OLS}} + c\|\beta\|^2.$$

Objective function.

$$Q(\beta).$$

regularization. (Penalty term)

; $\mathbf{X}^T \mathbf{X} = \mathbf{r}_{XX}$ and $\mathbf{X}^T \mathbf{Y} = \mathbf{r}_{YX}$,
 vinent of simple correlation between



★ Point of view (II): Minimization with constraints \Leftrightarrow P.O.V.I.

$$\underline{b}^R = \underset{\underline{B}}{\operatorname{argmin}} Q(\underline{B}) = \underset{\underline{B}}{\operatorname{argmin}} \|y - \underline{x}\underline{B}\|^2 \quad \text{subject to: } \|\underline{B}\|^2 \leq t.$$

★ Simultaneously minimizing residuals AND minimizing \underline{B} .

$$\begin{array}{l} \text{If } C=0: \underline{b}^R = \underline{b}^{\text{OLS}} \\ \text{If } C \rightarrow \infty: \underline{b}^R = \underline{0} \end{array} \quad \Rightarrow \text{As } C \uparrow, \text{ bias } \uparrow.$$

P.O.V.I = P.O.V.II
under KFT conditions.

★ P.O.V.III: Linear Algebra.

$$\text{M.C.} \Rightarrow \det(\underline{X}^T \underline{X}) \approx 0.$$

$$= \pi_1 \cdot \pi_2 \cdot \pi_3 \cdots \pi_K \approx 0.$$

\hookrightarrow at least one eigenvalue $\pi_i \approx 0$.

$$\text{Eigenvalues } (\underline{X}^T \underline{X} + C \underline{I}) \Rightarrow \text{eigenvalues } \pi_1 + C; \\ \pi_2 + C; \\ \vdots \\ \pi_K + C;$$

$$\det(\underline{X}^T \underline{X} + C \underline{I}) = (\pi_1 + C)(\pi_2 + C) \cdots (\pi_K + C) \neq 0.$$

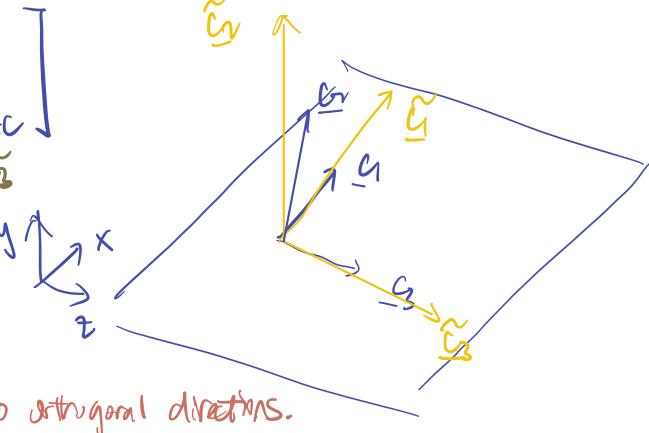
★ P.O.V.III: Geometric Interpretation.

$$\text{say } \underline{X}^T \underline{X} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad C_1 = C_2$$

$$\underline{X}^T \underline{X} + C \cdot \underline{I} = \begin{bmatrix} 1+C & 0 & 1 \\ 0 & 1+C & 0 \\ 1 & 0 & 1+C \end{bmatrix}$$

Add the ridge correction term makes these columns less linearly dependent.

Consider a general case where the C_j 's are linearly dependent.



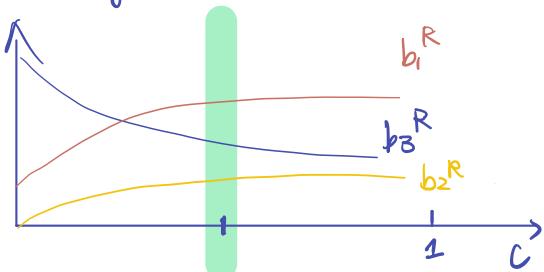
→ How to choose c ?

Plot b_k^R for different values of $c \in [0, 1]$

Plot VIF_k for different values of $c \in [0, 1]$

$VIF_k = \frac{1}{1 - R_k^2} \rightarrow$ regressing X_k onto the other predictors. (\uparrow if correlated).

Plot b_k against c .



Choose this c .

$b(c)$, b_k^R are "stable".

Downsides:

- Choice of c is subjective.
- No distributional results (CIs, PIs)

↳ Bootstrapping

(II) Robust Regression (Ch. 11.3)

↳ Useful when we have clusters of outliers.

OLS is sensitive to influential obsn's

Idea: create regression dampens influential obsn's

⇒ Robust regression I: LAR / LAD / L^1

least absolute residuals

least absolute deviations

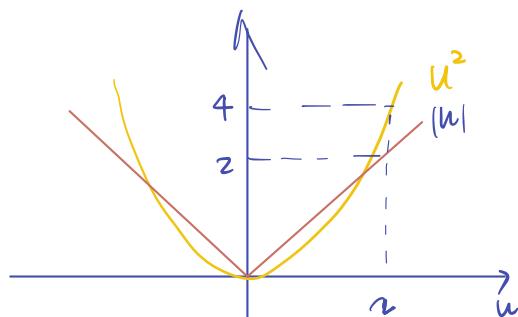
L^1 norm

$$\text{obj. } \min_{\beta} \sum_{i=1}^n |y_i - (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_{p-1} x_{ip})|$$

$$L^2 \text{ norm: } \|u\|_2 = \sqrt{\sum_{i=1}^n u_i^2} \leftarrow \text{OLS} = \|u\|_2$$

$$L^1 \text{ norm: } \|u\|_1 = \sum_{i=1}^n |u_i| \leftarrow \text{LAR} = \|u\|_1$$

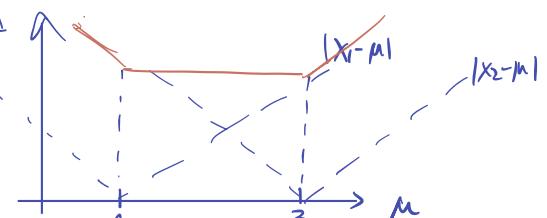
LAD / L^1



Downside:

not unique solutions.

Compare $f = T_i \cdot |x_i - \mu|$ & $f = T_i \cdot (x_i - \mu)^2$
let's consider $x_1 = 1, x_2 = 3$.



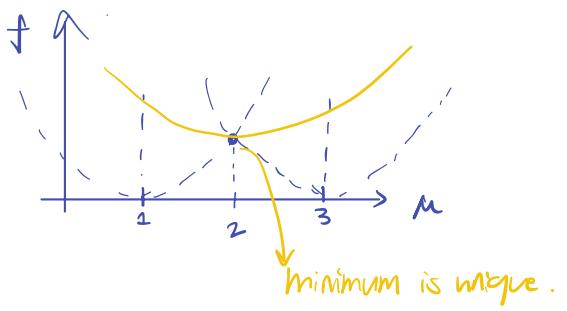
•) Robust Regression II: LMS - Least median of squares

$$\min_{\beta} \text{median} \left[y_i - (\beta_0 + \beta_1 x_1 + \dots + \beta_p x_{p-1}) \right]^2$$

$$\min_{\beta} \text{median} (\epsilon_i^2)$$

	mean	median
1 2 3	2	2
1 2 10	≈ 4	2

robust.



Downside: Hard to optimize. Median not differentiable.

•) Robust Regression III: IRLS. iteratively reweighted least squares.

Similar to WLS, but influential obs's are downweighted while others treated the same.

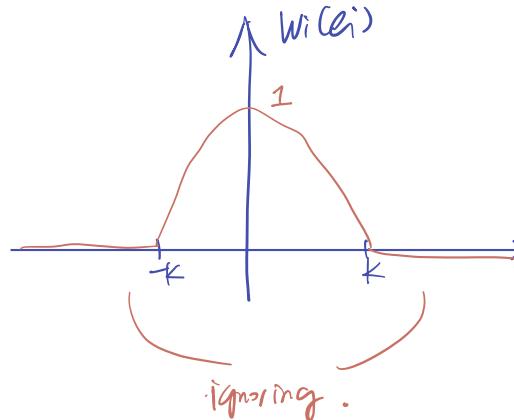
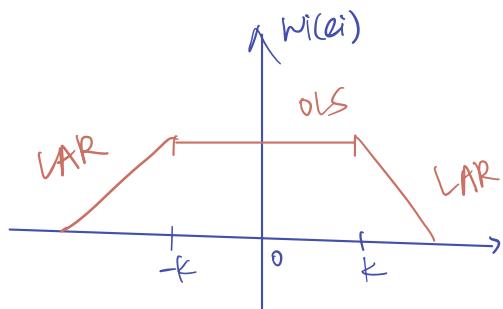
$$\min \sum_i w_i(\epsilon_i)^2$$

Huber

$$w_i(\epsilon_i) = \begin{cases} 1, & |\epsilon_i| < k \\ \frac{k}{|\epsilon_i|}, & |\epsilon_i| > k \end{cases}$$

Bisquare

$$w_i(\epsilon_i) = \begin{cases} \left[1 - \left(\frac{|\epsilon_i|}{k} \right)^2 \right]^2, & |\epsilon_i| < k \\ 0, & |\epsilon_i| > k. \end{cases}$$



* Lecture 20. 2023/11/13 (I).

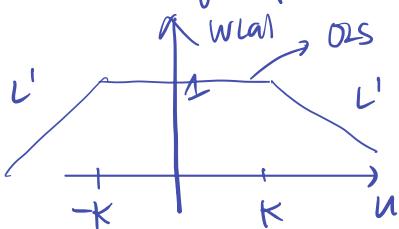
(I) Robust Regression - remedial measure for outliers (clusters)

LAD / LAR / L^1

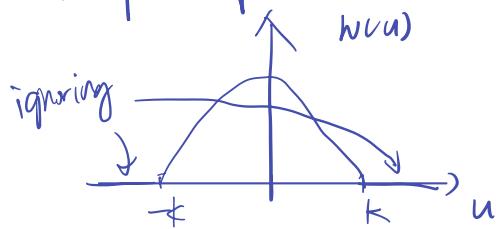
LMS (median)

IRLS

Huber weight func.



Bisquare function.



$$u_i = e_i / \text{MAD}$$

$\text{MAD} = \text{robust estimator of } \sigma$
= median absolute deviation.

$$= \frac{\sum_{i=1}^n |x_i - \text{median}|}{n}$$

* If $e_i \sim N(0, \sigma^2)$, $E(\text{MAD}) = \sigma$

Do IRLS iteratively:

- obtain initial weights with OLS

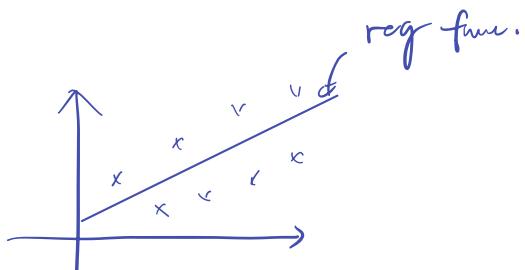
- iterate until b_k^{IRL} stabilize OR look at $w(u_i)$

$$\begin{matrix} e_i / u_i \\ \hat{y}_i \end{matrix}$$

Usually start w/ Huber, switch to bisquare later.

* Properties

- Requires knowledge of reg. func.
- If reg. func. unknown / difficult to establish
⇒ use nonparametric reg. methods.
(lowess, regression trees)
- Can identify multiple outliers.
- Rob. regression is used to confirm OLS: if $b_k^{IRL} \approx b_k^{OLS} \Rightarrow$ no influential obsns



Downside: can't obtain CIs, PI's, etc.

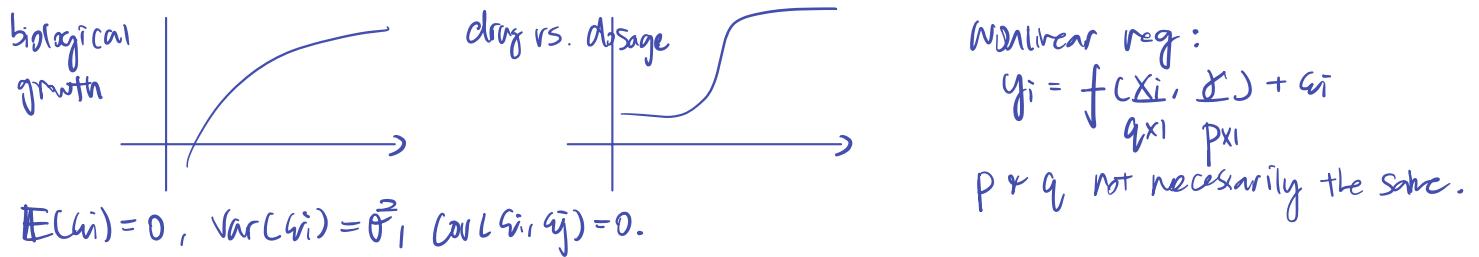
solution: bootstrap

we eliminated influence of y outlying obsn's

For X outlying obsn's: bounded influence regression

(II) Nonlinear regression. (Ch 13)

lots y_i 's are nonlinear func. of predictors.



Nonlinear reg:
 $y_i = f(x_i, \beta) + \epsilon_i$
 $g_{x_i} \quad p_{x_i}$

$p \neq q$, not necessarily the same.

Examples: Eq. 1: $f(x_i, \beta) = x_i^\top \beta \leftarrow \text{MLR}$

Eq. 2: exponential reg. model: (A) $y_i = \beta_0 \exp(\beta_1 x_i) + \epsilon_i \quad q=1, p=2$
(B) $y_i = \beta_0 + \beta_1 \exp(\beta_2 x_i) + \epsilon_i \quad q=1, p=3.$

Eq. 3: Logistic Reg.: $y_i = \frac{\beta_0}{1 + \beta_1 \exp(\beta_2 x_i)} + \epsilon_i \quad q=1, p=3.$

• Linear (MLR)

• Intrinsically linear (A). $f(x_i, \beta) = \beta_0 \exp(\beta_1 x_i)$

$$\log f = \underbrace{\log \beta_0}_{\beta_0} + \beta_1 x_i = \tilde{\beta}_0 + \beta_1 x_i$$

How to obtain β ? use OLS, MLE.

DLS: $Q(\beta) = \sum_{i=1}^n (y_i - f(x_i, \beta))^2$ (Normal Eqn.)
 $\hat{\beta} = \arg \min_{\beta} Q(\beta) \quad \frac{\partial Q}{\partial \beta_k} = 2 \sum (y_i - f(x_i, \beta)) \frac{\partial f}{\partial \beta_k} (x_i, \beta) = 0. \leftarrow$

MLE: $L(\beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum (y_i - f(x_i, \beta))^2\right)$

For MLR, e.g. $x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip-1} \end{bmatrix}$

(III) Probit / Logit Regression (Ch. 14)

Multilinear Regression when y is categorical / quantitative with 2 outcomes.

e.g. $y_i = \text{High BP / low BP.}$

= voted yes / no.

= bought stocks / didn't

= took a drug / didn't

$$y_i = \begin{cases} 1 & \text{w/ prob. } \pi_i \\ 0 & \text{w/ prob. } 1 - \pi_i \end{cases}$$

$$y_i \sim \text{Bernoulli}(\pi_i)$$

$$\text{If use SLR: } y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$E(y_i) = \beta_0 + \beta_1 x_i = 1(\pi_i) + 0(1-\pi_i) = \pi_i$$

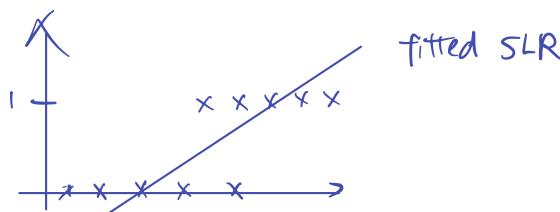
Has to be in $[0, 1]$

For all x_i depending on x_i , this does not happen.

e.g. $X = \text{alcohol use}$

$y^c = \text{duration of pregnancy}$ (continuous resp. variable)

$$\text{SLR: } y_i^c = \beta_0^c + \beta_1^c x_i + \epsilon_i^c, \quad \epsilon_i^c \sim N(0, \sigma^2_c)$$



Dichotomization: $y_i^c = \begin{cases} 1, & y_i^c \leq 38 \text{ weeks} \\ 0, & y_i^c > 38 \text{ weeks} \end{cases}$ preterm
full term.

$$\pi_i = P(y_i=1) = P(y_i^c \leq 38) = P(\beta_0^c + \beta_1^c x_i + \epsilon_i^c \leq 38) = P\left(\frac{\epsilon_i^c}{\sigma_c} \leq \frac{38 - \beta_0^c}{\sigma_c} - \frac{\beta_1^c}{\sigma_c} x_i\right)$$

$\epsilon_i^c \sim N(0, \sigma_c^2)$

$$z = \frac{\epsilon_i^c}{\sigma_i} \sim N(0, 1)$$

$$\pi_i = P(z \leq \beta_0^* + \beta_1^* x_i) = \Phi(\beta_0^* + \beta_1^* x_i)$$

$$\mathbb{E}(y_i) = \pi_i = \Phi(\beta_0^* + \beta_1^* x_i)$$

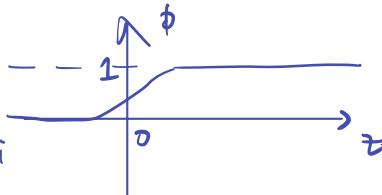
Non linear response func.: $f(x_i, \beta) = \Phi(\beta_0^* + \beta_1^* x_i)$

Probit mean func.

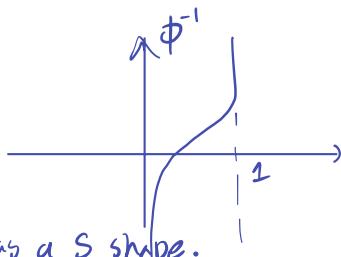
$$\pi_i = \Phi(\beta_0^* + \beta_1^* x_i)$$

$$\pi_i' = \Phi^{-1}(\pi_i) = \beta_0^* + \beta_1^* x_i$$

Probit probit transform



$$\phi^{-1}: [0, 1] \rightarrow (-\infty, \infty)$$



In general, $\phi(\beta_0^* + \beta_1^* x_i)$ has a S shape.

$$y_i' = 1 - y_i \Rightarrow \beta_0' = -\beta_0^* \\ \beta_1' = -\beta_1^*$$

If $\beta_1^* \uparrow$, ϕ^{-1} shifts right
If $\beta_1^* > 0$, ϕ^{-1} shifts up
If $\beta_1^* < 0$, ϕ^{-1} shifts down

* Logit Mean Response

Logistic R.V.: $f_L(w) = \frac{\exp(w)}{1 + \exp(w)} \quad \mathbb{E}(w) = 0, \text{Var}(w) = \frac{\pi^2}{3}$



Modify error term: $\epsilon_i \sim N(0, \sigma^2)$

$\epsilon_{L,i} \sim \text{logistic} (\text{mean}=0, \text{var}=\sigma_L^2)$

$$\mathbb{E} y_i = \pi_i = P(y_i^c = 1) = P(y_i^c \leq 38) = P(\epsilon_i^L \leq 38 - \beta_0 - \beta_1 x_i) = P\left(\frac{\epsilon_i^L}{\sigma} \leq \frac{38 - \beta_0 - \beta_1 x_i}{\sigma}\right)$$

$$= P(\epsilon_i^L \leq \beta_0^* + \beta_1^* x_i)$$

$\pi_i = F_L(\beta_0^* + \beta_1^* x_i)$ CDF of logistic R.V.

$$= F_L(\beta_0^* + \beta_1^* x_i)$$

$$F_L(w) = \frac{1}{1 + \exp(-w)} = \exp(w) / (1 + \exp(w)) \quad (\text{also S-shaped})$$

$$\pi_i = F_L(\beta_0^* + \beta_1^* x_i)$$

$$\pi_i' = F_L'(\pi_i) = \beta_0^* + \beta_1^* x_i$$

Lecture 21 2023年11月15日(三).

Logistic Regression

Probit mean response func.

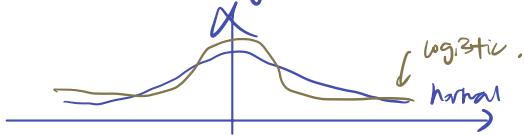
$$\begin{aligned} \mathbb{E}y_i &= \pi_i = \Phi(\beta_0 + \beta_1 x_i) \\ \hookrightarrow \text{squashes } \beta_0 + \beta_1 x_i \text{ into } [0, 1] \\ \pi_i' &= \Phi^{-1}(\pi_i) = \beta_0 + \beta_1 x_i \\ \hookrightarrow \text{stretches } (0, 1) \text{ into } [-\infty, +\infty] \end{aligned}$$

$y_i \sim \text{Bernoulli}(\pi_i)$

$$y_i = \begin{cases} 1 & \text{W.p. } \pi_i \\ 0 & \text{W.p. } (1 - \pi_i) \end{cases}$$

Logit mean response func.

errors ~ logistic R.V.'s



$$\begin{aligned} \mathbb{E}y_i &= \pi_i = F_L(\beta_0 + \beta_1 x_i) \\ \pi_i' &= F_L^{-1}(\pi_i) = \underline{\beta_0 + \beta_1 x_i} \end{aligned}$$

$$F_L(u) = \frac{1}{1 + \exp(-u)} = \text{cdf of logistic R.V.}$$

$$\pi_i' = F_L^{-1}(\pi_i) = \log \frac{\pi_i}{1 - \pi_i} \quad (\text{log odds})$$

$$\text{odds} = \frac{P(S)}{P(F)} = \frac{\pi_i}{1 - \pi_i}$$

$$\begin{aligned} \text{If } \pi_i = \frac{1}{2} \Rightarrow \text{odds} = 1 \\ \text{If } \pi_i = \frac{1}{10} \Rightarrow \text{odds} = \frac{1}{10} = \frac{1}{9} \end{aligned}$$

another option - for reference: error terms ~ Gumbel R.V.

$$\pi_i = \log(\log(1 - \pi_i)) = \beta_0 + \beta_1 x_i \quad (\text{3rd most common})$$

→ most common: Logistic regression.

→ Simple logistic regression.

$$\mathbb{E}y_i = \pi_i = \frac{\exp(\beta_0 + \beta_1 x_i)}{1 + \exp(\beta_0 + \beta_1 x_i)} \Leftrightarrow \pi_i' = \beta_0 + \beta_1 x_i$$

Coefficients are estimated via MLE:

$$P(y_i = 1) = \pi_i ; P(y_i = 0) = 1 - \pi_i$$

$$P(y_i = y) = \pi_i^y (1 - \pi_i)^{1-y}$$

$$\begin{aligned} L(y_1, \dots, y_n | \beta_0, \beta_1) &= \text{joint pmf of } y_1, \dots, y_n = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \\ &\quad \downarrow \text{given obsns.} \quad \downarrow \text{unknowns} \end{aligned}$$

↓ take the log

$$\ell(y_1, \dots, y_n | \beta_0, \beta_1) = \log L = \sum y_i \log \pi_i + (1 - y_i) \log (1 - \pi_i)$$

$$= \sum_{i=1}^n y_i \log \left(\frac{\pi_i}{1 - \pi_i} \right) + \log(1 - \pi_i)$$

$$\text{① log odds} = \pi_i^l = \beta_0 + \beta_1 x_i$$

$$= \sum_{i=1}^n \left\{ y_i (\beta_0 + \beta_1 x_i) - \log(1 + \exp(\beta_0 + \beta_1 x_i)) \right\}$$

→ Solve numerically to find maximizers.

$$\begin{aligned} b_0 &= \hat{\beta}_0^{\text{MLE}} \\ b_1 &= \hat{\beta}_1^{\text{MLE}} \end{aligned}$$

$$\hat{\pi}_i = \frac{\exp(b_0 + b_1 x_i)}{1 + \exp(b_0 + b_1 x_i)}$$

fitted value for the i th obsn

* Interpretation of b_1

$$\begin{aligned} \hat{\pi}_i^l(X_i) &= b_0 + b_1 x_i \\ \hat{\pi}_i^l(X_{i+1}) &= b_0 + b_1 x_i + b_1 \end{aligned} \Rightarrow$$

$$\text{odds} = \frac{\pi_i^l}{1-\pi_i^l}$$

$$\begin{aligned} b_1 &= \hat{\pi}_i^l(X_{i+1}) - \hat{\pi}_i^l(X_i) \\ b_1 &= \pi_i^l(2) - \pi_i^l(1) \\ &= \log(\text{odds}_2) - \log(\text{odds}_1) \\ &= \log \frac{\text{odds}_2}{\text{odds}_1} \end{aligned}$$

$$\text{odds ratio} = OR = \frac{\text{odds}(X_{i+1})}{\text{odds}(X_i)} = e^{b_1}$$

Increase of predictor value by 1 leads to increase of odds by $\exp(b_1)$.

$$b_1 = 0.1615 \Rightarrow \exp(b_1) = 1.175$$

As X increases by 1, odds of y are increased by 17.5%.

* Multiple Logistic Regression.

$$\mathbb{E}y_i = \pi_i = \frac{\exp(X^T \beta)}{1 + \exp(X^T \beta)} = (1 + \exp(-X^T \beta))^{-1}$$

$$\pi^l = \underbrace{\begin{pmatrix} X^T \beta \\ 1 \end{pmatrix}}_{n \times p, p+1 \text{ predictors}} \rightarrow \text{vector of predictors w/ intercept.}$$

$n \times p$, $p+1$ predictors

* Inference (Chp. 14.5)

We only have large sample results for $S^2(b_1)$ via theory of MLE.

$$S^2(b) \approx L(-g_{ij})_{p=k}^{-1} \text{ where } g_{ij} = \frac{\partial^2 \log L}{\partial \beta_i \partial \beta_j} \quad \text{Hessian of log-likelihood}$$

* Wald-Test (t -test)

$$H_0: \beta_k = 0 \text{ vs. } H_a: \beta_k \neq 0.$$

$$z^* = \frac{b_k}{S(b_k)}$$

CI for odds ratio: $OR = e^{\beta_k}$

$$\begin{aligned} |z^*| &= \left| \frac{b_k}{S(b_k)} \right| > z(1 - \frac{\alpha}{2}), \text{ rej. } H_0 \\ &\leq z(1 - \frac{\alpha}{2}), \text{ accept } H_0. \end{aligned}$$

$$\underbrace{e^{\beta_k \pm z(1 - \frac{\alpha}{2}) S(b_k)}}_{\text{CI for } \beta_k}$$

$$\hat{\pi}_h \pm \left(1 - \frac{\alpha}{2}\right) \sqrt{\hat{X}_h^T S^2(b) \hat{X}_h} = [L, U]$$

sd of mean response at X_h .

$$\text{For } \hat{\pi}_h : [L^*, U^*] \text{ where } L^* = (1 + e^{-L})^{-1} \quad U^* = (1 + e^{-U})^{-1}$$

$$\hat{\pi}_h = (1 + e^{-\hat{T}_h})^{-1}$$

$$X_h = [x_1 \ x_{h1} \ x_{h2} \ \dots \ x_{hp}]$$

* Polytomous Logistic Regression (Ch. 14.11, p608-611) Multinomial

As opposed to dichotomous (binomial)

Response variable y_i = qualitative / categorical with > 2 categories

$$\text{Eq. : } y_i = \begin{array}{l} \text{product A} \\ \text{OR} \\ \text{product B} \\ \text{OR} \\ \text{product C} \end{array} \quad \text{nominal}$$

nominal
ordinal - there is a natural ordering to classes.

$$y_i = \begin{array}{l} \text{mild illness} \\ \text{moderate illness} \\ \text{severe illness} \end{array} \quad \text{ordinal}$$

$$y_{ij} = \begin{cases} 1 & \text{if cat 1} \\ 0 & \text{o.w.} \end{cases}$$

$$y_{i2} = \begin{cases} 1 & \text{if cat 2} \\ 0 & \text{o.w.} \end{cases}$$

$$y_{i3} = \begin{cases} 1 & \text{if cat 3} \\ 0 & \text{o.w.} \end{cases}$$

$$\sum_{j=1}^3 y_{ij} = y_{i1} + y_{i2} + y_{i3} = 1$$

one hot encoding.

! Not dummy coding !

$$\pi_{ij} = P(y_{ij}=1) = \text{"prob. that cat } j \text{ is selected in } i \text{th response"}$$

One can show that:

$$\pi_{ij} = \frac{\exp(X_i^T b_j)}{1 + \sum_{k=1}^{J-1} \exp(X_i^T b_k)}$$

where $J = \#$ of categories y_i can assume.

we skipped category j .

Note: If $J=r \Rightarrow$ regular logistic regression.

Lecture 27 2013-11-27 (A-L)

Poisson Regression (Ch. 14.13, pp. 618-620)

y_i = count of "rare" events

y_i = # of trips family takes to mall in 1 month (1)

y_i = # of hospitalizations in a year (2)

x_i = # of members, income, distance, ... (1)

x_i = preexisting conditions, age, more (2)

$y \sim \text{Poisson}(\mu)$

P.m.f. $f(y) = \frac{\mu^y e^{-\mu}}{y!}$, $y=0, 1, 2, 3, \dots$ $\mathbb{E}y = \mu$, $\text{Var}(y) = \mu$.

Poisson (y_n) \approx Binomial (n, p) as $n \rightarrow \infty$, $p \rightarrow 0$.

$\mu = n-p$. "rare events"

Poisson regression

y_i are independent Poisson r.v.s with parameter μ_i (= expected value μ_i)

$$\mu_i = \mu(x_i, \beta)$$

\downarrow Vector of Predictors.

Vector of Predictors

(*) 3 options for $\mu(x_i, \beta)$:

$$\mu(x_i, \beta) = x_i^\top \beta$$

$$\mu(x_i, \beta) = \exp(x_i^\top \beta)$$

$$\mu(x_i, \beta) = \log(x_i^\top \beta)$$

β is estimated via MLE:

$$L(\beta) = \prod_{i=1}^n \frac{\mu(x_i, \beta)^{y_i} e^{-\mu(x_i, \beta)}}{y_i!}$$

$$l(\beta) = \log L(\beta) = \sum_i y_i \log(\mu(x_i, \beta)) - \sum_i \mu(x_i, \beta) - \sum_i y_i \log(y_i!)$$

\downarrow

maximize to find β !

$$\text{Fitted values } \hat{\mu}_i = \mu(x_i, \beta)$$

General Linear Models (GLMs) Ch. 14.15.

(bern., normal, poisson, exp,...)

1) y_1, \dots, y_n are resp. var's coming from exponential family of dist'n's

2) A linear predictor $x_i^\top \beta = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$

3) A link function $g \rightarrow$ monotonic & differentiable.

$$x_i^\top \beta = g(\mu_i) \quad \text{mean of } y_i, \quad \mathbb{E}(y_i) = \mu_i$$

$$\text{Var}(y_i) = \sigma_i^2 \quad (\text{nonconstant}) \quad \text{but } \sigma_i^2 = f(\mu_i) = f(g^{-1}(x_i^\top \beta))$$

Eq. \rightarrow Normal error model $g(\mu_i) = \mu_i$

$$x_i^\top \beta = g(\mu_i) = \mu_i \Rightarrow x_i^\top \beta = \mu_i \Rightarrow \mathbb{E}y_i = x_i^\top \beta$$

\rightarrow Logistic Regression: $g = F_L^{-1}(\mu_i) \rightarrow$ cdf of logistic R.V.

$$x_i^\top \beta = g(\mu_i) \Rightarrow x_i^\top \beta = F_L^{-1}(\mu_i) = F_L^{-1}(\pi_i)$$

$$\Rightarrow \mathbf{x}_i^T \boldsymbol{\beta} = \ln \mu_i = \frac{\log \mu_i}{1 - \mu_i}$$

•) Poisson Regression: $g(\mu_i) = \log(\mu_i)$
 $\mathbf{x}_i^T \boldsymbol{\beta} = \log(\mu_i) \Rightarrow \mu_i = \exp(\mathbf{x}_i^T \boldsymbol{\beta})$

* One factor ANOVA (Chp. 16.)

Factor: Predictor Variable = X qualitative / categorical

levels of factor = treatments = values of X . X_1, X_2, \dots, X_r

E.g. factor = dosage of drug.

① treatments = low, medium, high.

y = effectiveness of 3 dosages of drug to reduce blood pressure.

② factor = brand

treatments = B_1, B_2, B_3, B_4 .

y = absorbency of paper towel.

y_{ij} = response of the j th replicate of i th treatment

treat 1

treat 2

treat 3

$y_{11} \ y_{12} \ \dots \ y_{1n_1}$

$y_{21} \ y_{22} \ \dots \ y_{2n_2}$

$y_{31} \ y_{32} \ \dots \ y_{3n_3}$

n_i = # of replicates in treatment i

total # of observations = $N_T = \sum n_i$

* Cell means model for ANOVA (CMM)

$$y_{ij} = \mu_i + \epsilon_{ij}$$

respons. variable in j th replicate / trial for i th factor level or treatment.

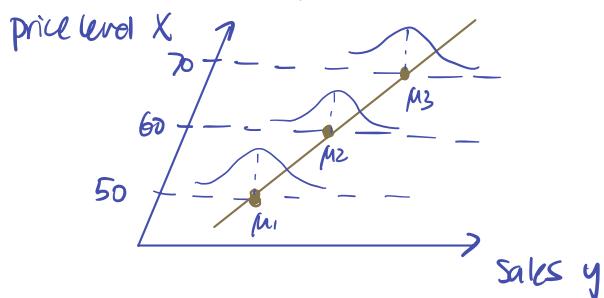
μ_i = parameters

$\epsilon_{ij} \sim \text{indep. } N(0, \sigma^2)$

$i = 1, \dots, r$ # of factor levels / treatments

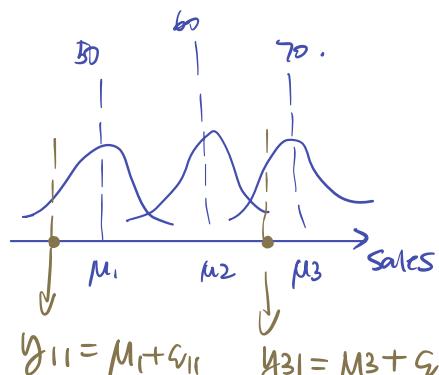
$j = 1, \dots, n_i$ # of observations in treatment i

$\sum n_i = N_T$ total # observations.



$y_{ij} \sim N(\mu_i, \sigma^2)$, independent.

CELL Means Model



* ANOVA as a linear model

treatments $r=3$, $n_1 = n_2 = n_3 = 2$

$$y = \begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \\ y_{31} \\ y_{32} \end{bmatrix} \quad \beta = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} \quad \Rightarrow y_{ij} = \mu_i + \epsilon_{ij} \Leftrightarrow y = X\beta + \epsilon.$$

$$X = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad \epsilon = \begin{bmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{31} \\ \epsilon_{32} \end{bmatrix}$$

$$\mathbb{E}y = X\beta = \begin{bmatrix} \mu_1 \\ \mu_1 \\ \mu_2 \\ \mu_2 \\ \mu_3 \\ \mu_3 \end{bmatrix}$$

μ_i = population mean for treatment i
 σ^2 = "variability" of responses in population.

$$\text{Var } y = \text{Var } \epsilon = \sigma^2 I_6 \hookrightarrow 6 \times 6 \text{ identity matrix.}$$

* What if we just did linear regression w/ dummy coding?

$$\begin{aligned} x_{ij1} &\left\{ \begin{array}{ll} 1, & \text{if } T_1 \\ 0, & \text{o.w.} \end{array} \right. \\ x_{ij2} &\left\{ \begin{array}{ll} 1, & \text{if } T_2 \\ 0, & \text{o.w.} \end{array} \right. \end{aligned} \quad \text{defining dummy coding}$$

$$y_{ij} = \beta_0 + \beta_1 x_{ij1} + \beta_2 x_{ij2} + \epsilon_{ij}$$

$$\mathbb{E}y_{j1} = \beta_0 + \beta_1 \quad \text{(if treatment 1) measures effect of treatment 1.}$$

$$\mathbb{E}y_{j2} = \beta_0 + \beta_2$$

Regression approach to ANOVA or effects model.

CMM \rightarrow cell means μ_1, μ_2, μ_3 .

We study only fixed effects ANOVA: factor levels are fixed / predetermined.

We do NOT study random effects ANOVA.

factor levels \rightarrow chosen randomly from some population.

Notation:

$$y_{fi} = \sum_{j=1}^{n_i} y_{ij} = \text{sum of observations in } i\text{th treatment}$$

$$\bar{y}_{fi} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} = \bar{y}_{i..} = \text{avg. of obs'n's in } i\text{th treatment.}$$

$$\bar{y}_{i..} = \frac{1}{n_i} y_{fi} = \frac{n_i}{\sum_{f=1}^r n_f} \bar{y}_{f..}$$

Overall mean $\bar{y}_{...} = \frac{1}{NT} \sum_{f=1}^r \sum_{j=1}^{n_f} y_{fj}$ \hookrightarrow weighted sum of treatment averages.

→ Estimate of cell means (LS, MLE). Lecture 13.

1) LS.

$$Q(\mu_1, \dots, \mu_r) = T_i T_j (y_{ij} - \mu_i)^2 \rightarrow \text{objective function.}$$

$$= T_j (y_{1j} - \mu_1)^2 + T_j (y_{2j} - \mu_2)^2 + \dots + T_j (y_{nj} - \mu_r)^2$$

minimize wrt cell means μ_i

$$\hat{M}_{iLS} = \bar{y}_i. \quad \text{fitted values } \hat{y}_{ij} = \bar{y}_i.$$

In vector notation, $\hat{y} = \underline{x} b$, where $b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$

M_1
 M_2
 M_3

2) MLF

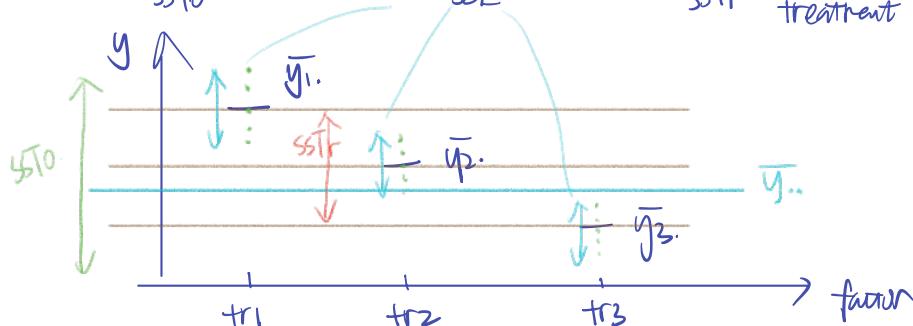
$$L(\mu_1, \dots, \mu_r, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \cdot \sum_{i,j} (y_{ij} - \mu_i)^2\right)$$

maximize $L \rightarrow \hat{m}_i^{\text{MLE}} = \hat{m}_i^{\text{LS}} = \bar{y}_i$. for $i=1, \dots, r$.

Residuals: $e_{ij} = y_{ij} - \hat{y}_{ij} = y_{ij} - \bar{y}_i$. $\sum_{j=1}^{n_i} e_{ij} = 0$. for $i=1, \dots, r$

In vector notation: $\mathbf{e} = \mathbf{y} - \mathbf{x}\mathbf{b}$.

$$\sum_{j=1}^{T_i} (y_{ij} - \bar{y}_{i..})^2 = \sum_{j=1}^{T_i} (y_{ij} - \bar{y}_j)^2 + \sum_{j=1}^{T_i} M_{ij} (\bar{y}_j - \bar{y}_{i..})^2$$



ANOVA Table :

SS	df	MS	$\text{E}(MS)$
SSE	$n_T - r$	$\text{MSE} = \text{SSE} / (n_T - r)$	σ^2
SSTr	$r-1$	$\text{MST}_r = \text{SST}_r / (r-1)$	$\sigma^2 + \bar{T}_i \frac{n_i (M_i - M_0)^2}{(r-1)}$
SSTo	$n_T - 1$		→ Where $M_0 = \bar{T}_i = \frac{\sum n_i}{n_T} M_i$ Weighted mean.

Scenario 1): $SSTr > SSE$ \rightarrow treatment works. i.e. M_i are different.

Scenario 2): $SSTr < SSE \rightarrow$ Treatment prob. doesn't work, i.e. μ_i are "Same"

$H_0: \mu_1 = \mu_2 = \dots = \mu_r$, i.e. treatments don't work.

H_a : at least one μ_i is different, so treatment works.

under H_0 , $SSE/\sigma^2 \sim \chi_{n-r}^2$

$SSTr/\sigma^2 \sim \chi_{r-1}^2$

$$\frac{SSTr}{\sigma^2/(r-1)} = \frac{MSTr}{MSF} \sim F(r-1, n-r) \quad F^* = \frac{MSTr}{MSF}$$

under H_0 .

If $F^* > F_{1-\alpha, r-1, n-r}$ conclude H_a

between group diff. ($SSTr$) \gg within group difference (SSE).

If $F^* \leq F_{1-\alpha, r-1, n-r}$ conclude H_0 .

Intuition: $E(MSE) = \sigma^2$
 $E(MSTr) = \sigma^2 + \frac{\sum_{i=1}^r n_i(\mu_i - \mu_0)^2}{r}$ ≈ 0 on average when we conclude H_0 .

under H_0 , $MSE \approx MSTr$ on average

under H_a , $MSE < MSTr$ on average

in M.L.R., F^* was MSR/MSE

* Factor effects model (Ch. 1b.7)

Recall, in cell means model $y_{ij} = \mu_i + \epsilon_{ij}$ Factor effect model:

$\mu_i = \mu_0 + \frac{(\mu_i - \mu_0)}{T_i} = \mu_0 + T_i$ ↑ i th factor level effect / i th treatment effect.

constant component of all obsn's

so $y_{ij} = \mu_0 + T_i + \epsilon_{ij}$, $i = 1, \dots, r$, $j = 1, \dots, n_i$

* Definition of μ_0 = unweighted mean.

$$\mu_0 = \frac{1}{r} \sum_{i=1}^r \mu_i \Rightarrow \sum_{i=1}^r T_i = 0$$

$$\text{e.g. } (r=3) : \frac{1}{3} (70 + 60 + 53) = 61$$

$$T_1 = 9, T_2 = -1, T_3 = -8 \\ \rightarrow T_1 + T_2 + T_3 = 0.$$

→ treatment 1 increases the mean of obsn's from that group by 9 units.

* weighted mean $\mu_0 = \sum_{i=1}^r w_i \mu_i$, $\sum_i w_i = 1$
 $\Rightarrow \sum_i w_i T_i = 0$.

Choose w_i : according to sample size: $w_i = n_i/n_T \rightarrow \sum_i n_i/n_T = \frac{1}{n_T} \sum_i n_i = n_T/n_T = 1$.

OR

according to the importance of the treatments to the experimenter.

In CMM,

$$H_0: \mu_1 = \mu_2 = \dots = \mu_r$$

H_a : not all μ_i are equal

In factor effects model,

$$H_0: \tau_1 = \tau_2 = \dots = \tau_r = 0.$$

v.

H_a : not all $\tau_i \neq 0$.

"treatments have no effect on response variable."

Lecture 24 2023/12/4(一).

One factor ANOVA

regression approach to ANOVA

multiple inferences for factor level means.

Regression Approach to ANOVA (Ch. 1b.8).

3 models : (A) Factor effects with unweighted mean.

(B) " " " weighted "

(C) Cell means model.

$$y = X\beta + \varepsilon$$

??

(A) Factor effects w/ unweighted mean

$$y_{ij} = \mu_0 + \tau_i + \varepsilon_{ij}, \quad i=1, \dots, r$$

common effect of treat i .

Need to estimate $(r+1)$ params. $\mu_0 \& \tau_i, i=1, \dots, r$

↳

only r are independent: $\tau_i \tau_j = 0$

$$\Rightarrow \tau_r = -\tau_1 - \tau_2 - \dots - \tau_{r-1}$$

e.g. $r=3$; $n_1=n_2=n_3=2$

$$y = \begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \\ y_{31} \\ y_{32} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \mu_0 \\ \tau_1 \\ \tau_2 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \end{bmatrix}$$

$$\varepsilon = \begin{bmatrix} \varepsilon_{111} \\ \varepsilon_{112} \\ \varepsilon_{211} \\ \varepsilon_{212} \\ \varepsilon_{311} \\ \varepsilon_{312} \end{bmatrix}$$

$$\mu_0 - \tau_1 - \tau_2 = \mu_0 + \tau_3$$

$$y = X\beta + \varepsilon \Leftrightarrow y_{ij} = \mu_0 + \tau_i + \varepsilon_{ij}$$

$$\hat{y} = X\hat{\beta} = \begin{bmatrix} \mu_0 + \tau_1 \\ \mu_0 + \tau_1 \\ \mu_0 + \tau_2 \\ \mu_0 + \tau_2 \\ \mu_0 - \tau_1 - \tau_2 \\ \mu_0 - \tau_1 - \tau_2 \end{bmatrix}$$

$$y_{ij} = \mu_0 + \tau_i x_{ij1} + \tau_r x_{ij2} + \dots + \tau_{r-1} x_{ijr-1} + \varepsilon_{ij}$$

Intercept

slope coeffs.

$$y_{rj} = \mu_0 - \tau_1 - \tau_2 - \dots - \tau_{r-1} + \varepsilon_{rj} = \mu_0 + \tau_r + \varepsilon_{rj}$$

$$X_{ij1} = \begin{cases} 1 & \text{if treat 1} \\ -1 & \text{if treat r} \\ 0 & \text{o.w.} \end{cases}$$

$$X_{ij(r-1)} = \begin{cases} 1 & \text{if treat r-1} \\ -1 & \text{if treat r} \\ 0 & \text{o.w.} \end{cases}$$

"effect" coding

L.S. estimates: $\hat{\mu}_0 = \frac{1}{n} \sum_{i=1}^r \bar{y}_{ji}$, $\hat{T}_i = \bar{y}_{ji} - \hat{\mu}_0$

$$\beta_1 \quad \beta_2 \quad \dots \quad \beta_{r-1}$$

$$\left\{ \begin{array}{l} H_0: T_1 = T_2 = \dots = T_{r-1} = 0 \\ H_a: \text{at least one } T_i \neq 0. \end{array} \right. \rightarrow \text{Can test this using the overall model FBS test.}$$

$$F^* = \frac{MSR}{MSE}$$

(B) Factor effects with weighted mean.

$$\mu_0 = \sum_i w_i \bar{y}_{ji} = \sum_i \frac{n_i}{n_T} \bar{y}_{ji}, \quad n_T = \sum_i n_i$$

$$\sum_{i=1}^r \frac{n_i}{n_T} T_i = 0 \Rightarrow T_r = -\frac{n_1}{n_T} T_1 - \frac{n_2}{n_T} T_2 - \dots - \frac{n_{r-1}}{n_T} T_{r-1}$$

↳ weighted by sample size.

$$Y_{ij} = \mu_0 + T_1 X_{ij1} + T_2 X_{ij2} + \dots + T_{r-1} X_{ij(r-1)}$$

$$X_{ij1} = \begin{cases} 1 & \text{if treat 1} \\ -\frac{n_1}{n_T} & \text{if treat r} \\ 0 & \text{o.w.} \end{cases}$$

$$X_{ij(r-1)} = \begin{cases} 1 & \text{if treat (r-1)} \\ -\frac{n_{r-1}}{n_T} & \text{if treat r} \\ 0 & \text{o.w.} \end{cases}$$

Note: if $n_1 = n_2 = \dots = n_r$, then we get (A).

e.g. if $r=4$, $n_1 = n_2 = n_3 = 5$, $n_4 = 4$.

$$X_{ij1} = \begin{cases} 1 & \text{if treat 1} \\ -\frac{1}{5} & \text{if treat 4} \\ 0 & \text{o.w.} \end{cases}, \quad X_{ij2} = \begin{cases} 1 & \text{if treat 2} \\ -\frac{1}{5} & \text{if treat 4} \\ 0 & \text{o.w.} \end{cases}, \quad X_{ij3} = \begin{cases} 1 & \text{if treat 3} \\ -\frac{1}{4} & \text{if treat 4} \\ 0 & \text{o.w.} \end{cases}$$

$H_0: T_1 = T_2 = \dots = T_{r-1} = 0$.

$H_a: \text{at least one } T_i \neq 0$.

$$F^* = \frac{MSR}{MSE}$$

(C) Cell means model.

$$Y_{ij} = M_i + C_{ij} \Leftrightarrow Y = X\beta + \varepsilon$$

$$\text{Where } \mathbf{1}\mathbf{1}_{n_1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \mathbf{D}_{n_1} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \mathbf{N}_1 \mathbf{X}_0.$$

$$\beta = \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_r \end{bmatrix}, \quad X = \begin{bmatrix} \mathbf{1}\mathbf{1}_{n_1} & \mathbf{D}_{n_1} & \cdots & \mathbf{D}_{n_1} \\ \mathbf{0}_{n_2} & \mathbf{1}\mathbf{1}_{n_2} & \mathbf{D}_{n_2} & \cdots \mathbf{D}_{n_2} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{n_r} & \cdots & \mathbf{D}_{n_r} & \mathbf{1}\mathbf{1}_{n_r} \end{bmatrix}$$

Eq.: Consider $r=3$ & n_1, n_2, n_3

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \begin{bmatrix} \bar{y}_{1\cdot} \\ \bar{y}_{2\cdot} \\ \bar{y}_{3\cdot} \end{bmatrix} \quad X_i = \begin{cases} 1 & \text{if treat 1} \\ 0 & \text{o.w.} \end{cases} \quad X_r = \begin{cases} 1 & \text{if treat r} \\ 0 & \text{o.w.} \end{cases}$$

no Intercept!

$$Y_{ij} = \mu_1 X_{ij1} + \mu_2 X_{ij2} + \dots + \mu_r X_{ijr} + \epsilon_{ij} \quad (\star\star)$$

$$H_0: \mu_1 = \mu_2 = \dots = \mu_r$$

$H_a:$ at least one μ_i is \neq than the others.

NOT the same as the overall model fit test.

use general linear test.

$$H_0: \text{reduced model } Y_{ij} = \mu_c + \epsilon_{ij}.$$

$$H_a: \text{full model } (\star\star)$$

$$\text{Reduced: } y = \mathbf{x}\beta + \epsilon, \quad \mathbf{x} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \beta = [\mu_c]$$

$$F^* = \frac{\frac{SSE(R) - SSE(F)}{SSE(F)}}{\frac{df(F)}{df(R) - df(F)}} \quad \begin{array}{l} \text{df}(F) \leftarrow n_T - r \\ \text{df}(R) - \text{df}(F) \leftarrow n_T - 1 - (n_T - r) = r - 1 \end{array}$$

$$= \frac{MS_{Tr}}{MSE}$$

★ Inference for factor level means

$$\text{CMM: } Y_{ij} = \mu_i + \epsilon_{ij}$$

Before, we could conduct: $H_0: \mu_1 = \mu_2 = \dots = \mu_r$
vs.

$H_a:$ not all μ_i are $=$.

Which μ_i if μ_j ?

OR What's the CI for μ_i ?

OR What's the CI for $(\mu_i - \mu_j)$?

★ Estimation / testing for single factor level mean μ_i

$$\text{L.S.: } \hat{\mu}_i = \bar{y}_{i\cdot} \quad Y_{ij} \sim N(\mu_i, \sigma^2), \text{ indep.}$$

$$\bar{y}_{i\cdot} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} \Rightarrow E(\bar{y}_{i\cdot}) = \mu_i \quad \text{Var}(\bar{y}_{i\cdot}) = \frac{\sigma^2}{n_i} \quad S^2(\bar{y}_{i\cdot}) = \frac{MSE}{n_i}$$

$$\bar{y}_{i\cdot} \sim N(\mu_i, \sigma^2/n_i) \quad (\bar{y}_{i\cdot} - \mu_i) / S(\bar{y}_{i\cdot}) \sim t(n_T - r)$$

(1) we use full sample n_T to estimate σ^2 .

d.f.s = $n_T - r$
because $SSE/\sigma^2 \sim \chi^2_{n_T - r}$
 $SSE = \sum_{i,j} (y_{ij} - \bar{y}_{i\cdot})^2$ $\hookrightarrow r$

$\star \star$ CI_(1-α) For $\mu_i = \bar{y}_{i\cdot} \pm t(1-\alpha/2, n_T-r) S(\bar{y}_{i\cdot})$, $i=1, \dots, r$

Hypothesis Test: $H_0: \mu_i = c$, vs. $H_a: \mu_i \neq c$.

$$t^* = \frac{\bar{y}_{i\cdot} - c}{S(\bar{y}_{i\cdot})} \sim t(n_T-r)$$

If $|t^*| \leq t(1-\alpha/2, n_T-r)$ $\Rightarrow H_0$
 " > " $\Rightarrow H_a$.

* Estimation / Testing for difference of 2 factor level means.

Define $D = \mu_i - \mu_j \Rightarrow \hat{D} = \bar{y}_{i\cdot} - \bar{y}_{j\cdot}$

$$\mathbb{E}(\hat{D}) = \mathbb{E}(\bar{y}_{i\cdot}) - \mathbb{E}(\bar{y}_{j\cdot}) = \mu_i - \mu_j = D$$

$$\text{Var}(\hat{D}) = \text{Var}(\bar{y}_{i\cdot}) + \text{Var}(\bar{y}_{j\cdot}) - \underbrace{2\text{Cov}(\bar{y}_{i\cdot}, \bar{y}_{j\cdot})}_{= 0} = \sigma^2 \left(\frac{1}{n_i} + \frac{1}{n_j} \right)$$

$$\hat{s}^2(\hat{D}) = \text{MSE}\left(\frac{1}{n_i} + \frac{1}{n_j}\right) = 0 \quad \text{because obs'n's in } i\text{th cell } \perp\!\!\!\perp \text{ obs'n's in } j\text{th cell.}$$

$$\Rightarrow \frac{\hat{D} - D}{S(\hat{D})} \sim t(n_T-r)$$

CI_(1-α) For $D = \hat{D} \pm t(1-\alpha/2, n_T-r) S(\hat{D})$

$$\hat{D} = \bar{y}_{i\cdot} - \bar{y}_{j\cdot}$$

$$S(\hat{D}) = \sqrt{\text{MSE}\left(\frac{1}{n_i} + \frac{1}{n_j}\right)}$$

Hypothesis Test: $H_0: D=0$, vs. $H_a: D \neq 0$

$$(\mu_i = \mu_j) \quad (\mu_i \neq \mu_j)$$

$$t^* = \frac{\hat{D}}{S(\hat{D})} \sim t(n_T-r)$$

Lecture 15 12/06 (E)

Multiple comparisons

Say $r=3$. Test 1: $H_0: \mu_1 = \mu_2$ vs. $H_a: \mu_1 \neq \mu_2$.

Test 2: $H_0: \mu_1 = \mu_3$ vs. $H_a: \mu_1 \neq \mu_3$.

Test 3: $H_0: \mu_2 = \mu_3$ vs. $H_a: \mu_2 \neq \mu_3$.

Each test conducted at $\alpha = 0.05$, FTR - "failing to reject"

$$P(\text{FTR } H_0 \text{ in } T_1 / H_0 \text{ is true}) = 1 - \alpha = 0.95$$

$$P(\text{FTR } H_0 \text{ in } T_1 \text{ & } T_2 \text{ & } T_3 / H_0 \text{ is true}) \approx 0.95^3 = 0.857$$

$$P(\text{rejecting at least once in } T_i's / H_0) = 1 - 0.857 = 0.143 = 14\%$$

Type I error for 3 tests = "family-wise" error (FWE).

$\Rightarrow FWE \leq \alpha$ (this is what we want)

→ Tukey Multiple Comparison Procedure. (= "Honestly Significant Differences") = HSV in R.

Defn: Studentized Range Distribution (SRD).

Consider i.i.d. sample $y_1, \dots, y_r \sim N(\mu, \sigma^2)$

$$W = \max_i(y_i) - \min_i(y_i) = \text{range.}$$

Estimate σ^2 by s^2 with $v \rightarrow n$

Define new r.v.: $q(r, v) = W/s$ where $r = \# \text{ obsns in samples}$

$$v = \# \text{ d.f.s of } s^2$$

$q(r, v) = \text{studentized range.}$

Percentiles of $q(r, v)$ are in Table B.9 (KNNL)

$$\text{e.g. } P\left(\frac{W}{s} = q(5, 10) \leq 4.65\right) = 0.95$$

$$\hookrightarrow \text{in R. } q(0.95, 5, 10)$$

$$1 - \alpha \quad r \quad v$$

Simultaneous CI for all possible differences.

$$\hat{D} \pm T \cdot S(\hat{D}), \quad \hat{D} = \bar{y}_i - \bar{y}_j, \quad S(\hat{D}) = \sqrt{\text{MSE}\left(\frac{1}{n_i} + \frac{1}{n_j}\right)}, \quad T = \frac{1}{\sqrt{v}} q(1-\alpha, r, n_T - r)$$

$\hookrightarrow \frac{r(r-1)}{2}$ CIs if we have r factor levels.

$$\frac{\hat{D}}{S(\hat{D})} = \frac{\bar{y}_i - \bar{y}_j}{S(\hat{D})} \leq \frac{\max_i(\bar{y}_i) - \min_i(\bar{y}_i)}{S(\hat{D})} \quad \text{worst case scenario.}$$

H.T.: (Simultaneous): $H_0: \mu_i = \mu_j$ vs. $H_a: \mu_i \neq \mu_j$ for all pairs.

$$q^* = \frac{T \cdot \hat{D}}{S(\hat{D})}$$

$$|q^*| \leq q(1-\alpha, r, n_T - r), \quad H_0.$$

Overall type I error $\leq \alpha$.

Bonferroni CI $1-\alpha$ (Simultaneous)

$$\hat{D} \pm B \cdot S(\hat{D}), \quad B = t(1-\alpha/2g, n_T - r), \quad g = \# \text{ comparisons.}$$

$$\text{If all pairwise comparisons: } g = \frac{r(r-1)}{2}$$

→ Tukey is better, less conservative.

If $q < \frac{k(k-1)}{2}$, then Bonf. could be better than Tukey.

★★ Two factor ANOVA ("Two way ANOVA")

Case where y depends on 2 categorical variables.

- o) $y = \text{Sales}$
- o) Pred. Var 1 = factor A, e.g. Price level (\$50, \$60, \$70) # levels = a.
- o) Pred. Var 2 = factor B, e.g. Marketing (TG, FB) campaign # levels = b.
total number of treatments: ab.

We shouldn't do one-factor ANOVA: we fail to assess the joint effects of both factors A & B.

		fact B	$i=1$	$j=2$
		fact A		
$i=1$		y_{111}, y_{112}, \dots	y_{121}, y_{122}, \dots	
$i=2$		y_{211}, y_{212}, \dots	y_{221}, y_{222}, \dots	
$i=3$		y_{311}, y_{312}, \dots	y_{321}, y_{322}, \dots	

$$\text{CMM: } y_{ijk} = \mu_{ij} + \alpha_{ijk}$$

$$\text{Effects: } y_{ijk} = \mu_{..} + \underline{\alpha_i} + \underline{\beta_j} + (\underline{\alpha\beta})_{ij} + \varepsilon_{ijk}.$$

Main effect of
factor A at level i

Main effect of
factor B at level j

new term c
(not the multiplication) interaction of A
at level i & B at level j