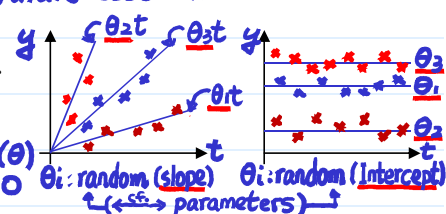


Generalized Least Square (GLS) → Recall OLS estimator is BLUE under Gauss-Markov condition

- model: $Y = X\beta + \epsilon$, $E(\epsilon) = 0$ and $\text{var}(\epsilon) = \sigma^2 I \Rightarrow \epsilon$: uncorrelated and constant variance
no normality $\text{Var}(Y) = \sigma^2 I$ $[Y]$ Q: how to model correlation in sampling model? (X, Y) vs. Y|X (LNp.5-7)
- Q: what if $\text{var}(\epsilon) \neq \sigma^2 I$? ϵ may have non-constant variance and/or are correlated, e.g., correlation between observations in Y. $[\text{Var}(\epsilon) = \text{Var}(Y)]$

- time series correlation [e.g., $\epsilon_t \sim \text{ARMA}(r, m)$] ① data observed over time
Auto-regression ② previous observations correlated with future observation
- growth curve model, repeated measurement model [e.g., several observations taken from same person, or same unit] data y_i, y_j from same unit $\text{cov}(y_i, y_j) = \text{cov}(\theta_i + \delta_i, \theta_j + \delta_j) = t_i t_j \text{Var}(\theta)$

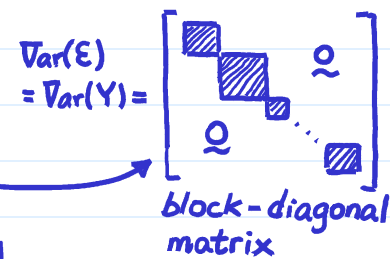


- spatial correlation [e.g., data taken over contiguous geographical areas: census tracts, countries, or states in a country. Nearby areas are often much alike] cf. nearby time points

split plot design

- nested errors [e.g., M sets of observations, each set from common production run, from same/common equipment, or from same survey-taker]

random effect model in DOE



- Consider the case $\text{var}(\epsilon) = \sigma^2 \Sigma$, where $\Sigma (\neq I)$ is known but σ^2 is unknown, i.e., we know the correlation and relative variance between the errors but we don't know the absolute scale

off-diagonal part of Σ diagonal part of Σ

Σ is a covariance matrix ($\because \text{Var}(\epsilon/\sigma) = \Sigma$) D: orthogonal matrix $DD^T = I$
 Because $\Sigma_{n \times n}$ is symmetric and positive definite, we can write $\Sigma = SS^T$, where S is an $n \times n$ nonsingular matrix (by Cholesky or spectral decompositions) $\Sigma^{1/2}(\Sigma^{1/2})^T$ not unique

Note: If $\text{Var}(Y) = \sigma^2 \Sigma$, we can find a matrix A s.t. $\text{Var}(AY) = E^*(A Y Y^T A^T) = A E^*(Y Y^T) A^T = (A \Sigma A^T) \sigma^2 = \sigma^2 I$

$Y = X\beta + \epsilon \Rightarrow S^{-1}Y = S^{-1}X\beta + S^{-1}\epsilon \Rightarrow Y' = X'\beta + \epsilon'$ where $Y' = S^{-1}Y, X' = S^{-1}X, \epsilon' = S^{-1}\epsilon$ and still a linear model (known response & predictors)

$E(\epsilon') = 0$ and $\text{var}(\epsilon') = \text{var}(S^{-1}\epsilon) = S^{-1}\text{var}(\epsilon)S^{-T} = S^{-1}\sigma^2 SS^T S^{-T} = \sigma^2 I$

can do OLS on Y' & X'

- GLS: find $\hat{\beta}$ that minimize Gauss-Markov conditions $(\Sigma^{1/2})^T(\Sigma^{-1/2})$

$\epsilon'^T \epsilon' = (Y' - X'\beta)^T (Y' - X'\beta) = (Y - X\beta)^T S^{-T} S^{-1} (Y - X\beta) = (Y - X\beta)^T \Sigma^{-1} (Y - X\beta)$

$0 = \hat{\epsilon}'^T \hat{Y}' = \hat{\epsilon}'^T \Sigma^{-1} \hat{Y}'$

$\hat{\beta} = (X'^T X')^{-1} X'^T Y' = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$
 $\hat{\beta} \equiv \hat{\beta}_{GLS}$ for original data X, Y

$\text{var}(\hat{\beta}) = \sigma^2 (X'^T X')^{-1} = \sigma^2 (X^T \Sigma^{-1} X)^{-1}$

Note1: $\epsilon'^T \epsilon'$, $\hat{\beta}$, and $\text{var}(\hat{\beta})$ are invariant to the choice of S.
 Q: What depends on the choice of S?
 Ans: $\hat{Y}' = S^{-1} \hat{Y}, \hat{\epsilon}' = S^{-1} \hat{\epsilon}$
 Note2: $\hat{\beta} = \hat{\beta}_{OLS}$ if $\Omega[\Sigma^{-1} X] = \Omega[X]$

GLS is like OLS regressing

$Y' = S^{-1}Y$ on $X' = S^{-1}X$ $\epsilon' = S^{-1}\epsilon$ $\sigma^2 I$ $\sigma^2 \Sigma$

- Q: why should not use ordinary least square when $\text{var}(\epsilon) \neq \sigma^2 I$?

model: $Y = X\beta + \epsilon$, $E(\epsilon) = 0$ and $\text{var}(\epsilon) = \sigma^2 \Sigma$, OLS estimator $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y$

$E(\hat{\beta}_{OLS}) = \beta$, $\text{var}(\hat{\beta}_{OLS}) = \sigma^2 (X^T X)^{-1} X^T \Sigma X (X^T X)^{-1} \Rightarrow$ unbiased but variance not minimized (Note: $c^T \hat{\beta}$ is BLUE of $c^T \beta \Rightarrow \text{var}(c^T \hat{\beta}) \leq \text{var}(c^T \hat{\beta}_{OLS})$)

$\hat{\beta}_{GLS}$ apply Gauss-Markov Thm on $Y' = X' \beta + \epsilon'$

- diagnostics (residual analysis) should be applied on $Y' - X' \hat{\beta} = S^{-1} (Y - X \hat{\beta}) = S^{-1} \hat{\epsilon} \leftrightarrow \epsilon$ because ϵ' are i.i.d. but not ϵ

- The practical problem is that Σ may not be known. It's usually necessary to make some assumptions and examine the residuals to estimate Σ (check lab for an example, IRWLS)

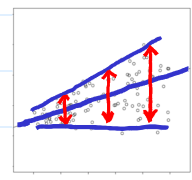
❖ Reading: Faraway (2005, 1st ed.), 6.1 ❖ Further reading: D&S, 9.2

add one more assumption than GLS
 $\text{Var}(\epsilon) = \sigma^2 \Sigma \neq \sigma^2 I$

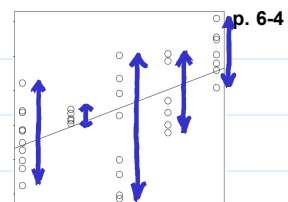
Weighted Least Square (WLS) *know relative variance*

- Sometimes, the errors are uncorrelated, but have unequal variance where the form of the inequality is known ($\Rightarrow \Sigma$ is diagonal, it's a special case of GLS), example:

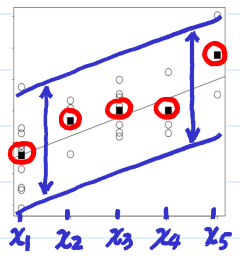
- error variance proportional to a function of predictors
 [e.g., $\text{var}(\epsilon_i) = x_i^2 \sigma^2 (\Rightarrow w_i = 1/x_i^2)$ or $\text{Var}(y_i) \propto [E(y_i)]^2$ or $\propto E(y_i) \leftarrow$ Poisson y_i 's]



- data with replicates, which show a pattern of unequal variance [e.g., $\text{var}(\epsilon_i) \approx$ sample variance of observations with same $x_i (\Rightarrow w_i = 1/(\text{sample variance}))$]



- the observed y_i 's are actually averages of several observations. [e.g., suppose y_i is the average of n_i observations, $\text{var}(\epsilon_i) = \sigma^2/n_i (\Rightarrow w_i = n_i)$]
- ① at x_i , observe y_{ij} , $j=1, \dots, n_i$ ② $y_{ij} = X\beta_i + \epsilon_{ij}$, $\epsilon_{ij} \text{ iid } N(0, \sigma^2)$
- ③ offer summarized data $(x_i, \bar{y}_i) \Rightarrow \bar{y}_i$'s are of constant variance $\Rightarrow \text{Var}(\bar{y}_i) = \sigma^2/n_i \leftarrow \bar{y}_i$'s not constant variance



- ϵ : uncorrelated, but not constant variance $\Rightarrow \Sigma$ is diagonal. Write

$\text{Var}(\epsilon) = \sigma^2 \Sigma$ parameter known

$$\Sigma = \begin{pmatrix} 1/w_1 & 0 & \dots & 0 \\ 0 & 1/w_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1/w_n \end{pmatrix} \Rightarrow \Sigma^{-1} = \begin{pmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w_n \end{pmatrix}$$

where w_i 's ($\propto 1/\text{var}(\epsilon_i)$) are called weights \leftarrow 權重

relative variance = $\frac{1}{w_1} \cdot \frac{1}{w_2} \cdot \dots \cdot \frac{1}{w_n}$

low weight \Leftrightarrow high variance; high weight \Leftrightarrow low variance

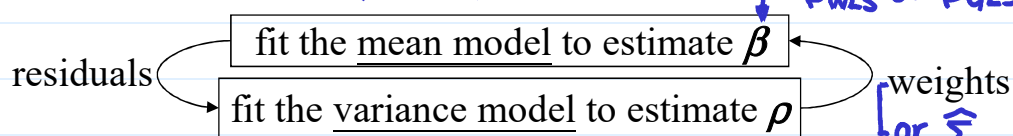
- $S^{-1} = \text{diag}(\sqrt{w_1}, \dots, \sqrt{w_n})$, then $\underline{S} = S S^T$ $\hat{\beta}_{WLS} = (X^T X)^{-1} X^T Y = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$
- \Rightarrow OLS regress $\underline{S}^{-1} Y$ (i.e., $\sqrt{w_i} y_i$) on $\underline{S}^{-1} X$ ($\sqrt{w_i} x_i$) (Note. the column of ones, i.e., intercept needs to be replaced with $\sqrt{w_i}$)
- \Rightarrow convenient for regression package without a weighted options

$RSS(X', Y') = \sum_i w_i [y_i - (X\beta)]^2$ (LNp.6-2)

a linear combination of the components in $Y' = [\sqrt{w_i} y_i]$

- Q: Why observations with smaller variance should be multiplied by heavier weight? intuitive interpretation? *pts with small variance have more influence on $\hat{\beta}_{WLS}$ pull the fitted line toward these data points*
- iteratively re-weighted least squares (IRWLS): In all the previous examples, weights (or Σ in GLS) are assumed known. Q: what if $\text{var}(\epsilon_i)$ is not completely known, what weights should we use? Q: where can you find the information of weights? *Ans: residuals or Σ*

procedure to estimate β in GLM



- Example: $\text{var}(\epsilon_i) = \rho_0 + \rho_1 x_{i1}$
- 1. start with $w_i = 1$

$E(\epsilon_i) = 0$
 $\text{Var}(\epsilon_i) = E(\epsilon_i^2) = \rho_0 + \rho_1 x_{i1}$

quasi-likelihood in GLM

- 2. use weighted least square to estimate β
- 3. use the residuals to estimate ρ_0 and ρ_1 , perhaps by regressing residuals² on x_1
- 4. re-compute the weights and go to 2. Continue until convergence

Garroll & Ruppert (1988)

Problems: converge? how is the inference about β affected? d.f.=? ...etc

alternative approach: jointly estimate the mean and variance parameters using likelihood based method (in R, use `glm()` function in the `nlme` library)

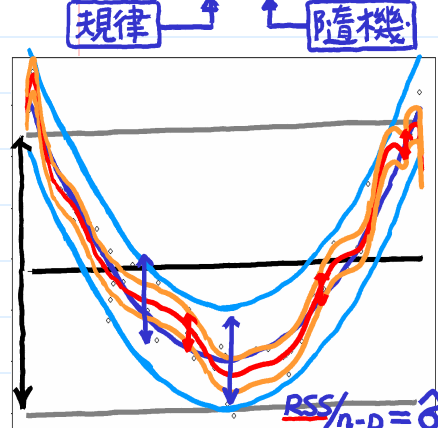
Reading: Faraway (2005, 1st ed.), 6.2 Further reading: D&S, 9.2

Testing for Lack of Fit

fitted model

Recall goodness of fit: R^2 & $\hat{\sigma}$ ($Y \neq \hat{Y}$)
 Note: good fit \rightarrow not necessary a right X

- model: $Y = X\beta + \epsilon$, Q: many choices of X , how can we tell the chosen X fits the data?



Q: what is "fit"? what "fit" is appropriate? example:

- data generated from the "true model" --- a 2nd-order polynomial of x $E(y|x) = \beta_0 + \beta_1 x + \beta_2 x^2$
- solid line (fitted model=1st-order): not capture the quadratic pattern in data *規律 \rightarrow 隨機 \rightarrow might produce biased $\hat{\beta}$*
- dashed line (fitted model=2nd-order): OK
- dotted line (fitted model=8th-order): fluctuated, fitted values and data are too close *規律 \leftarrow 隨機*

Why not compare R^2 to true R^2 ?

R^2 increases depends on the range of X (LN p.3-18)

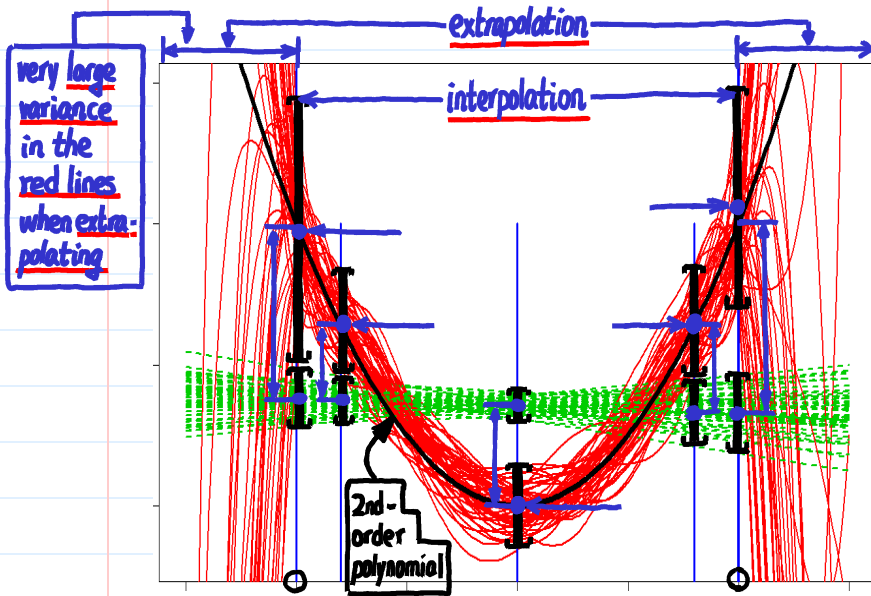
- X too simple \Rightarrow not enough to explain the mean structure in data \Rightarrow lack of fit (under-fitting) \Rightarrow $\hat{\sigma}^2$ over-estimated

- X too complex \Rightarrow will explain the variation caused by errors, in addition to the mean structure \Rightarrow overfit the data \Rightarrow $\hat{\sigma}^2$ under-estimated
- too simple & too complex could coexist in a fitted model*
- $Y = X^* \beta^* + \epsilon^*$
 $= X\beta + \epsilon = X\hat{\beta} + \hat{\epsilon}$
- problem of overfitting
 [current data] [future data] \rightarrow cross validation

Q: what statistic carry the information about lack of fit or overfit? Ans: $\hat{\sigma}^2$

- Repeat the simulation 50 times $\hat{Y} = P_{\Omega} Y = P_{\Omega} (X\beta + \epsilon)$
 - Black line: true model, 2nd-order polynomial
 - Green lines: fitted lines under under-fitting model, 1st-order polynomials
 - Red lines: fitted lines under over-fitting model, 8th-order polynomials

$span\{x_1\} = \Omega_1$
 $span\{x_1, x_2\} = \Omega_2^* = \Omega_1 \oplus \Omega_2$
 $span\{x_1, x_2, x_3\} = \Omega_3^* = \Omega_1 \oplus \Omega_2 \oplus \Omega_3$



under-fitting: $Y = X_1\beta_1 + \epsilon_1$
 correct: $Y = X_1\beta_1 + X_2\beta_2 + \epsilon_2$
 over-fitting: $Y = X_1\beta_1 + X_2\beta_2 + X_3\beta_3 + \epsilon_3$

Define $\Omega_1 = span\{x_1\}$, $\Omega_2 = \Omega_1^\perp$, $\Omega_3 = \Omega_2^\perp$
 $\Omega_2^* = span\{x_1, x_2\}$, $\Omega_3 = \Omega_2^* \cap \Omega_1^\perp$
 $\Omega_3^* = span\{x_1, x_2, x_3\}$, $\Omega_3 = \Omega_3^* \cap \Omega_2^{\perp}$

Define $\hat{Y}_1 = P_{\Omega_1} Y$, $\hat{Y}_2 = P_{\Omega_2} Y$, $\hat{Y}_3 = P_{\Omega_3} Y$

under-fitting: $\hat{Y} = \hat{Y}_1$
 correct: $\hat{Y} = \hat{Y}_1 + \hat{Y}_2$
 over-fitting: $\hat{Y} = \hat{Y}_1 + \hat{Y}_2 + \hat{Y}_3$

- Comparison
 - Under-fitting lines: small variance, large bias
 - Over-fitting lines: large variance, small bias
 - Interpolation vs. extrapolation

Recall.
 $MSE = \text{variance} + \text{bias}^2$

- Note: high R^2 (a measure of "goodness of fit") doesn't imply the model is a good fit (check lab for examples of high R^2 but lack of fit) \rightarrow overfitting
- graphical checking (based on residuals \leftarrow informal) (check a lab example)

a possible testing procedure for lack of fit (\Rightarrow compare $\hat{\sigma}^2$ to σ^2)

when σ^2 is known

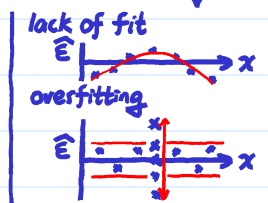
examples:

- σ^2 is known from past experience
- source of variation in ϵ is only measurement error, measuring device well understood, knowledge of measurement error inherent in an instrument or by definition

variance estimate under a fitted model: $Y = X\beta + \epsilon$

$Var(\epsilon^*)$
true error

check $\hat{\sigma}^2 \neq \sigma^2$



$z_i \rightarrow (y_i, \dots, y_{n_i}), \text{ obtain } (z_i, \bar{y}_i)$
 $\hat{\sigma}_i^2 \Rightarrow Var(\bar{y}_i / \frac{1}{n_i}) \approx 1/z_i$
 $Z = X\beta + \epsilon \Rightarrow Var(\epsilon) = I$

check LN p. 6-11

- each \bar{y}_i is an average of a very large number of observations, $\hat{\sigma}_i^2$ effectively estimated (see an example in lab)

test for lack of fit: fitted model: $Y = X\beta + \epsilon$, $\epsilon \sim N(0, \sigma^2 I)$, σ^2 : known

$H_0: E(Y) = X\beta$ is correct against $H_1: E(Y) = X\beta$ is too simple

- $(n-p)\hat{\sigma}^2 \sim \sigma^2 \chi_{n-p}^2$ if model is correct (under H_0), where $n = \#$ of observations, $p = \#$ of parameters in β \rightarrow null distribution

RSS

- test statistic: $(n-p)\hat{\sigma}^2 / \sigma^2 = RSS / \sigma^2$, compared with χ_{n-p}^2

exam over-fitting?

- reject if $(n-p)\hat{\sigma}^2 / \sigma^2 > \chi_{n-p}^2(1-\alpha)$ \leftarrow known \leftarrow Why? Why not 2-sided test?

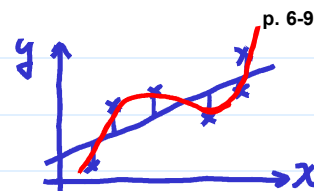
- if a lack of fit is found, a new (more complex) model is needed

when σ^2 is unknown (\Rightarrow can we use data to estimate it?)

- need "model-free method" to estimate σ^2 (i.e., free of the $E(Y)=X\beta$ assumption)

information is in residuals

residuals of what model?



Why? $\hat{E} = Y - X\hat{\beta}$ depends on the assignment of X

- because we want to use estimated σ^2 to justify whether $E(Y)=X\beta$ is suitable, the estimated σ^2 should have no relationship with the choice of X

- denote the estimated σ^2 under model-free method by $\hat{\sigma}_{p.e.}^2$, where p.e. stands for "pure error"

can treat the y 's on each distinct x_i as a one-sample case

- usually, only possible for data with replication (Q: why?)

Recall: one-sample case
 $Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu, \sigma^2)$
 $\nabla \epsilon \rightarrow \mu$
 $S^2 = \frac{1}{m-1} \sum_{j=1}^m (Y_j - \bar{Y})^2 \rightarrow \sigma^2$

- treat each group of data as one sample
- obtain $\hat{\sigma}^2$ of each sample
- pool the information in $\hat{\sigma}^2$'s together

how to estimate σ^2 (model-free)? Recall $\hat{\sigma}^2 = RSS/df$

$$SS_{p.e.} = \sum_{\text{distinct } x} \sum_{\text{within an } x} (y_{x,i} - \bar{y}_x)^2$$

$n - (\# \text{ of distinct } x_i)$

d.f. of p.e. = $\sum_{\text{distinct } x} (\# \text{ of replications} - 1)$

e.g. $(2-1) + (3-1) + (5-1) + (1-1) = 7$

$$\hat{\sigma}_{p.e.}^2 = \frac{SS_{p.e.}}{df_{p.e.}}$$

Why?

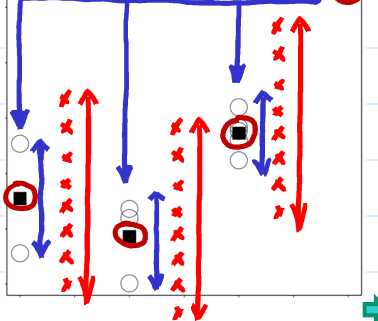
$$\hat{\sigma}_{pool}^2 = \frac{(m_1-1)\hat{\sigma}_1^2 + \dots + (m_k-1)\hat{\sigma}_k^2}{(m_1-1) + \dots + (m_k-1)}$$

test statistic:

($n = \#$ of observations, $p = \#$ of parameters in β)

RSS calculated from the model $E(Y)=X\beta$

of distinct x_i 's



←

$E(Y)=X\beta$	d.f.	SS	MS	F
Residual (ω)	$n-p$	RSS_{ω}	$RSS_{\omega}/n-p = \hat{\sigma}_{\omega}^2$	
Lack of fit	$n-p-df_{p.e.}$ ($k-p$)	$RSS - SS_{p.e.}$	$(RSS - SS_{p.e.})/(n-p-df_{p.e.})$ $(RSS_{\omega} - RSS_{\Omega})/k-p$	ratio of MS (compared to $F_{n-p-df_{p.e.}, df_{p.e.}}$)
Pure error (Ω)	$df_{p.e.}$ ($n-k$)	$SS_{p.e.}$ (RSS_{Ω})	$SS_{p.e.}/df_{p.e.} = \hat{\sigma}_{p.e.}^2$	null dist.

Recall: General form of F-test in LNp. 4-10



- Note: $\hat{\sigma}_{p.e.}^2$ is the estimate of σ^2 when we fit a saturated model to the data

k-sample model

$\Omega: E(Y_{ij}|x_i) = \mu_i, i=1, \dots, k$ # of distinct x_i 's
 $\mu = (\mu_1, \dots, \mu_k) \in \mathbb{R}^k, j=1, \dots, n_i (n_1 + \dots + n_k = n)$

$$Y = X\beta + \epsilon, X' = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}, \beta = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_k \end{bmatrix} \Rightarrow \hat{\sigma}^2 = \hat{\sigma}_{p.e.}^2$$

cf. LNp. 5-11 ($n=p$)

$\omega: E(Y_{ij}|x_i) = X\beta, \hat{Y}=?$

reduce the dimension of μ from k to p

of parameters = $p < k$

- alternative view:

Why cannot this test apply to no replicate data?

this is a comparison between the model of interest (i.e., $\omega: X\beta$) and a saturated model (Ω , whose R^2 reaches the maximum) that assigns a parameter to each unique combination of the predictors \Rightarrow standard F-testing for $H_0: \omega$ v.s. $H_1: \Omega \setminus \omega$

$\omega = \text{span}(X)$

not one if there are replicates (why? $\hat{E} \neq \rho$)

- need replication to make the test, but it's rare in obs'nal data
- possible solution: grouping (could be questionable \Rightarrow different grouping schemes may cause different conclusions)

