# Normality assumption

*Gauss-Markov Thm does not need this*

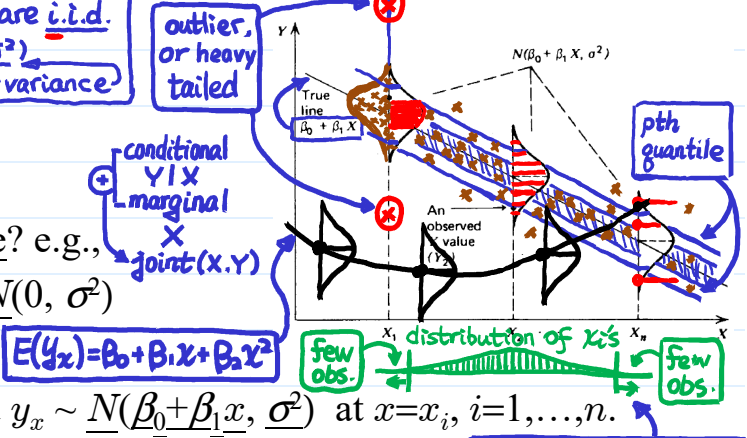- **Note**: up till <u>now</u>, **haven't assumed** any <u>distributional form</u> for $\varepsilon$. If we want to perform any <u>hypothesis tests</u> or make any <u>confidence intervals</u>, we will <u>need to do this</u>. The <u>usual assumption</u> is:

  *i.e. $\varepsilon_i$'s are i.i.d. $\sim N(0, \sigma^2)$ constant variance*

  | multivariate normal | $\varepsilon \sim N(0, \sigma^2 I)$ $\begin{bmatrix} \sigma^2 & 0 \\ 0 & \ddots \\ & & \sigma^2 \end{bmatrix}$ | *outlier, or heavy tailed* |

  ➤ <u>model</u>: $Y = X\beta + \varepsilon$, $\varepsilon \sim N(0, \sigma^2 I)$

  | matrix form | $Y \sim N(X\beta, \sigma^2 I)$ |

  *conditional $Y|X$ marginal $X$ joint $(X,Y)$*

  - **Q**: <u>what does</u> the <u>model</u> <u>describe</u>? e.g.,

  | functional form | $y_x = \beta_0 + \beta_1 x + \varepsilon_x$, $\varepsilon_x$'s $\sim$ i.i.d. $N(0, \sigma^2)$ |

  $\Rightarrow E(y_x) = \beta_0 + \beta_1 x$　　　$E(y_x) = \beta_0 + \beta_1 x + \beta_2 x^2$

  $\Rightarrow y_x$'s are <u>independent</u> and $y_x \sim N(\beta_0 + \beta_1 x, \sigma^2)$ at $x = x_i$, $i = 1, \ldots, n$.

  - **Q**: <u>how</u> should the <u>data</u> generated <u>from the model</u> <u>look like</u>?　*Recall $R^2$ in LN p.3-18*

  - **Q**: <u>when would it be appropriate</u> to <u>impose the inference</u> based on this <u>regression model</u> on the <u>underlying true model</u>? <u>Can we use it</u> when these exist <u>clear differences</u> between the <u>two models</u>, e.g., what if $y$ is a <u>discrete quantitative measurement</u>? *← check LN p. 2-3, $Y_i$'s "approximately" continuous.*

  **Ans**: yes when the <u>pdf shape</u> of the regression model can "<u>well approximate</u>" the <u>pdf/pmf/cdf shape</u> of true model. (<u>Key</u> is <u>how similar</u> the <u>2 models</u>)?

  > George Box: "<u>all models are wrong</u>, but <u>some are useful</u>"

---

➤ **Q**: <u>why Normal</u>?

$\varepsilon_i = \delta_1 + \delta_2 + \cdots + \delta_{i_k}$
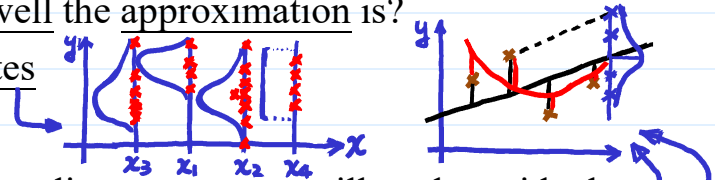
- <u>CLT</u> $\Rightarrow$ when <u>random error</u> is a <u>sum of many small random disturbances</u>

- <u>bell shape curve</u> is <u>common</u>

- from the viewpoint of <u>approximation</u>

- <u>good</u> <u>mathematical/statistical</u> properties

  *Note. Many distributions can be approximated by Normal, e.g., $Bin(m, p) \xrightarrow{d}$ Normal, as $m \to \infty$ $Poisson(\lambda) \xrightarrow{d}$ Normal, as $\lambda \to \infty$*

➤ **Q**: <u>how to examine</u> whether Normality assumption is <u>reasonable/suitable</u> for your data, in other words, <u>how well</u> the <u>approximation</u> is?

  - when you have <u>pure replicates</u>

  - when you have <u>no/few pure replicates</u>, you can still study <u>residuals</u>. However, the <u>validity</u> of the <u>study</u> is then based on <u>several assumptions</u>. Under the <u>circumstance</u>, what <u>rationale</u> can support the <u>use of Normality</u>?

➤ **Q**: <u>under what conditions</u>, the Normality assumption is <u>inappropriate</u>? *Check zero-inflated data*

  | GLM | · <u>qualitative response</u> → *counts* |
  | transfor-mation | · <u>quantitative discrete response</u> with only <u>few possible outcomes</u> |
  | | · <u>skewed error</u> |
  | | · <u>heavy tail error</u> |

  $\varepsilon_i$'s: *Normal / skewed*

  *e.g., response $y_i$'s have upper/lower bound*

  *e.g. $Bin(m,p)$ with small/large $p$.*

  *$\varepsilon_i$'s: Normal, heavy tail*

  *Normal → using (weighted) average of $y_i$'s heavy tail → causing problem in averaging*

  *OLS estimator is still valid*

＊ Some properties of (multivariate) Normal distribution　　p. 4-3

$E^*[(AZ+C)(AZ+C)^T]$
$= E^*[(AZ)(AZ)^T] = E^*[AZZ^TA^T]$

(N1). linear transformation of Normal is still Normal

$\boxed{\substack{n\times 1 \\ \text{vector}}}$ $Z \sim N(\mu, \Sigma)$ $\boxed{\substack{n\times n \\ \text{matrix}}}$ $\Rightarrow$ $AZ+c \sim N(A\mu+c, A\Sigma A^T)$ $\text{cov}(Z) = \begin{bmatrix} \text{cov}(Z_1) & \text{cov}(Z_1, Z_2) \\ \hline & \text{cov}(Z_2) \end{bmatrix}$

(N2). when $1^{\text{st}}$ and $2^{\text{nd}}$ moments are given, the Normal distribution is specified
　　　　　　　　　　　　　↳ i.e., mean vector & variance-covariance matrix

(N3). $Z = \begin{bmatrix} Z_1 \\ \hline Z_2 \end{bmatrix}$: Normal and uncorrelated (cov$(Z_1, Z_2)=0$) $\Rightarrow$ $Z_1, Z_2$ independent

$\boxed{\text{By (N3)}}$ $\underset{\boxed{\text{Normal}}}{W} = \begin{bmatrix} W_1 \\ W_2 \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} Z$, $\text{cov}(W_1, W_2) = E^*(W_1 W_2^T)$　　　$\boxed{\substack{\text{can be} \\ \text{generalized} \\ \text{to } k}}$
　　　　　　　　　　　　　　　　　　　$= E^*(A_1 ZZ^T A_2^T) = A_1 E^*(ZZ^T) A_2^T$

$\boxed{\text{By (N3) \& (N4)}}$ (N4). $Z \sim N(\mu, \Sigma)$, $W_1 = A_1 Z$, $W_2 = A_2 Z$ $\Rightarrow$ $W_1, W_2$ are independent iff $A_1 \Sigma A_2^T = 0$
　　　　　　can be generalized to $k$ ↩　　If $\Sigma = \sigma^2 I$, then $A_1 \Sigma A_2^T = 0 \Leftrightarrow A_1 A_2^T = 0$ ↩

$\boxed{\substack{\text{length}^2 \\ \text{of } W_i}}$ (N5). $Z \sim N(\mu, \Sigma)$, $W_1 = A_1 Z$, $W_2 = A_2 Z,\ldots$, $W_k = A_k Z$, and $\text{cov}(W_i, W_j)=0$ for $i \neq j$
　　　$\Rightarrow$ $W_1^T W_1$, $W_2^T W_2,\ldots$, $W_k^T W_k$ are mutually independent ← $\boxed{\substack{\text{useful for the indepen-} \\ \text{dence between SS's.}}}$
　　　　　　　　　　　　　　　　Recall If $A_i$'s are projection matrices. →

(N6). $Z$: an $n\times 1$ random vector and $Z \sim N(\mu, \Sigma)$, then　　　　$\boxed{\substack{\Sigma^{-\frac{1}{2}}(Z-\mu) \sim N(0, I) \\ \text{↳ standardization } \boxed{\text{By (N1)}}}}$

$\boxed{(\Sigma^{-\frac{1}{2}})^T (\Sigma^{-\frac{1}{2}})}$ ▪ $(Z-\mu)^T \Sigma^{-1}(Z-\mu) \sim \chi_n^2$ if $\Sigma$ is non-singular ←

　　　　　▪ $(Z-\mu)^T \Sigma^-(Z-\mu) \sim \chi_r^2$ if $\Sigma$ is singular and has rank $r$ $(< n)$,
$\boxed{\text{not unique}}$　　　where $\Sigma^-$ is a generalized inverse of $\Sigma$, (i.e., $\Sigma \Sigma^- \Sigma = \Sigma$)

• Some properties of linear models when $\varepsilon \sim N(0, \sigma^2 I)$ :　　$\boxed{\substack{\text{The possible vectors of } Z \text{ only} \\ \text{occupy a } r\text{-dim space of} \\ \text{the } n\text{-dim space } \mathbb{R}^n.}}$

$\boxed{\text{By (N1)}}$ ➢ distribution of $Y$ $[=X\beta + \varepsilon] \sim N(X\beta, \sigma^2 I)$

　　　➢ distribution of $\hat{\beta}$ $[=(X^TX)^{-1}X^T Y] \sim N(\beta, (X^TX)^{-1}\sigma^2)$ ↳ e.g., $\hat{Y}, \hat{\varepsilon}$

---

$\boxed{\text{By (N1)}}$ distribution of $\hat{\varepsilon}$ $[=(I-H)Y=(I-H)\varepsilon] \sim N(0, (I-H)\sigma^2)$, which has a singular　　p. 4-4
covariance matrix $I-H$ with rank $n-p$ (Note: $\dim(\hat{\varepsilon})=n-p$)　　$\boxed{\substack{\text{Note. eigenvalues of} \\ I-H \subset \substack{n-p\ 1's \\ p\ 0's}}}$

$\boxed{\text{By (N6)}}$ distribution of $RSS$ $[=(n-p)\hat{\sigma}^2 = \hat{\varepsilon}^T\hat{\varepsilon} = \varepsilon^T(I-H)\varepsilon] \sim \sigma^2 \chi^2_{n-p}$
　　　　　　　　　　　$\boxed{(I-H)}$ ← $(I-H)^2 = I-H \Rightarrow (I-H)^- = I-H$

$\boxed{\text{By (N1)}}$ distribution of $\hat{Y}$ $[= X\hat{\beta} = HY] \sim N(X\beta, H\sigma^2)$, which has a singular covariance
matrix with rank $p$ (Note: $\dim(\hat{Y})=p$)　　$\boxed{\substack{\text{Note. eigenvalues of} \\ H \subset \substack{p\ 1's \\ n-p\ 0's}}}$

$\boxed{\text{By (N4)}}$ $\hat{\beta}$ is independent of $\hat{\sigma}^2$ (Note: $\text{cov}((X^TX)^{-1}X^T Y, (I-H)Y)=0$) $\overset{\hat{\varepsilon}}{\frown}$
$\boxed{A_1 A_2^T = 0}$ $\boxed{\hat{\varepsilon}}$ $\underset{A_1}{\frown}$ $\underset{A_2}{\frown}$ $A_1 \Sigma A_2^T = (X^TX)^{-1}X^T \sigma^2 I(I-H)$

$\boxed{\text{By (N4)}}$ $\hat{Y}$ is independent of $\hat{\varepsilon}$ (Note: $\text{cov}(HY, (I-H)Y)=0$) $\overset{\hat{\varepsilon}}{\frown}$ $= \sigma^2[(X^TX)^{-1}X^T(I-H)]$
↳cf. different from $\hat{Y}^T\hat{\varepsilon}=0$ (LN p.3-4) ← $\boxed{A_1^T A_2 = 0}$ $A_1$ ⌐$A_2$ $\neq$ $= 0$

　　➢ distribution of prediction for a new set of predictors, $x_0 = (g_1(x_{10},\ldots,x_{m0}), \ldots,$
$\boxed{\substack{E(y_{x_0}) \\ = x_0^T \beta}}$ $g_p(x_{10},\ldots,x_{m0}))^T$　　model: $y = \sum_{j=1}^{p} \beta_j \cdot g_j(x_1,\ldots,x_m) + \epsilon$ $\boxed{\substack{\text{data} \\ \text{matrix}}}$
↳parameter　　　cf. $\substack{\text{fitted model : } \boxed{\hat{\beta}_j} \\ }$　　$x_0^T\hat{\beta}$　　$\boxed{x_0: \text{model matrix}}$

$\boxed{\substack{\text{random} \\ \text{variable}}}$ ▪ mean response v.s. future observation (Q: what different?)

$\boxed{y_{x_0} = x_0^T \beta + \varepsilon}$ ▫ Example: average yield when $x=x_0$? and tomorrow's yield when $x=x_0$?
　　　▫ same predicted value $x_0^T\hat{\beta}$, but different distributions

$\boxed{\substack{\text{future} \\ \text{error}}}$ ◉ distribution of prediction error for mean response at $x_0$

$\boxed{\text{By (N1)}}$ $x_0^T\hat{\beta} - x_0^T\beta \sim N(0, (x_0^T(X^TX)^{-1}x_0)\sigma^2)$ $\underset{\text{sample size is large}}{\overset{\text{close to 0 when}}{\longrightarrow}}$

$\boxed{\substack{E(\varepsilon)=0 \\ \text{Var}(\varepsilon)=\sigma^2}}$ ◉ distribution of prediction error for future observations at $x_0$

　　$x_0^T\hat{\beta} - (x_0^T\beta + \varepsilon) \overset{\text{indep.}}{\sim} N(0, (x_0^T(X^TX)^{-1}x_0 + 1)\sigma^2)$

❖ **Further reading**: Seber (1977), *Linear Regression Analysis*, chapter 2, 3.4, 5.2, 5.3　$\underset{\text{sample size is large}}{\overset{\text{not reduced when}}{\longrightarrow}}$

➤ Example: bivariate normal $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N\left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \right)$  $\rho: cor(x_1, x_2)$
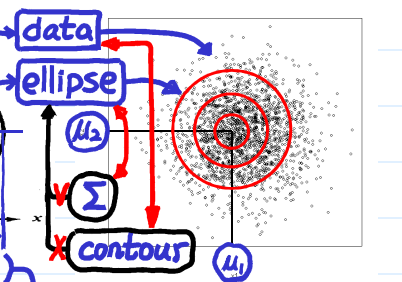
$\underbrace{\quad}_{\underset{\sim}{\mu}} \quad \underbrace{\quad}_{\Sigma}$

(a) $\sigma_1 = \sigma_2$, $\rho = 0$ $\Rightarrow$ independent and equal variance

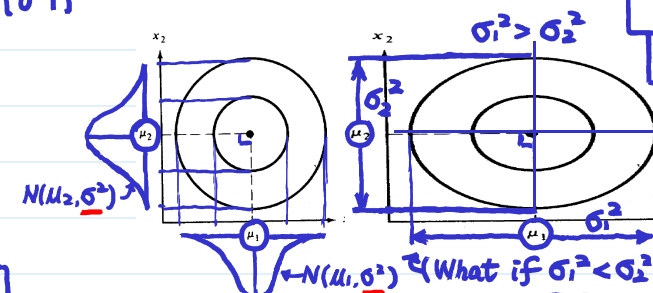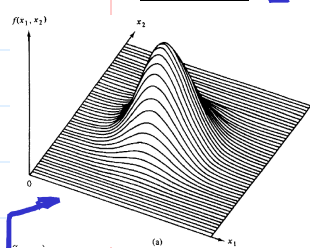joint pdf  $\Sigma = \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$        contour lines of the pdf        data generated from the pdf

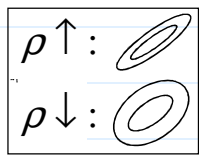**What if data not from normal?**



$\propto \exp\left[ -(\underset{\sim}{x} - \underset{\sim}{\mu})^T \Sigma^{-1} (\underset{\sim}{x} - \underset{\sim}{\mu}) \right]$

$N(\mu_2, \sigma^2)$      $\sigma_1^2 > \sigma_2^2$      $\sigma_2^2$      $\sigma_1^2$

$\leftarrow N(\mu_1, \sigma^2)$   **What if $\sigma_1^2 < \sigma_2^2$?**

**same marginal distributions**

data
ellipse
$\Sigma$
contour

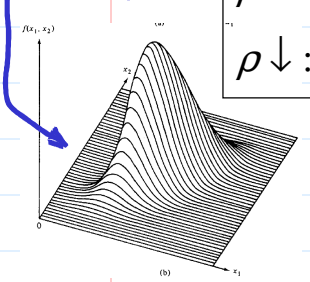**Q: how should the contour lines look like if $\sigma_1 \neq \sigma_2$?**

(b) $\sigma_1 = \sigma_2$, $\rho = 0.75$ $\Rightarrow$ correlated and equal variance

$\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$    $\rho \uparrow:$ ⬭    $\rho \downarrow:$ ⬭

$\sigma_1^2 > \sigma_2^2$

$N(\mu_2, \sigma^2)$
$x_1 = x_2$
$x_1 = -x_2$
$N(\mu_1, \sigma^2)$

**not parallel to $x_1 = x_2$ & $x_1 = -x_2$ any more**
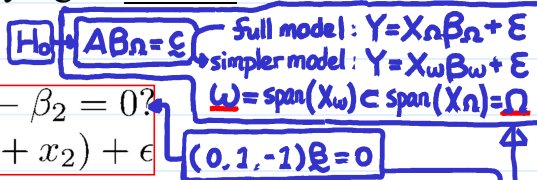
**quadratic form**

when $\sigma_1 = \sigma_2$, $\rho \neq 0$, the major/minor axis of the ellipse is parallel to $x_1 = x_2$ or $x_1 = -x_2$

contour of Normal pdf is an ellipse because it can be expressed as $(x - \mu)^T \Sigma^{-1} (x - \mu) = c$

---

$\beta \in \mathbb{R}^3$, $\Omega \subset \mathbb{R}^n$ but $dim(\Omega) = 3$ **hypothesis testings** (for $\beta$)

• **Q:** What questions is a hypothesis testing about $\beta$ trying to answer?

examples:  full model: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$   $H_0: A\beta_\Omega = c$

**full model: $Y = X_\Omega \beta_\Omega + \epsilon$**
**simpler model: $Y = X_\omega \beta_\omega + \epsilon$**
$\omega = span(X_\omega) \subset span(X_\Omega) = \Omega$

Q1: $\beta_1 = 0$? $(0,1,0)\beta = 0$
$\Rightarrow y = \beta_0 + \beta_2 x_2 + \epsilon$

Q2: $\beta_1 = \beta_2$? i.e., $\beta_1 - \beta_2 = 0$?
$\Rightarrow y = \beta_0 + \beta_1(x_1 + x_2) + \epsilon$   $(0,1,-1)\beta = 0$

**Ans**: Are all predictors needed? Can a simpler model still "well describe" the data?

• **Q:** Why a simpler model is preferred?

$\beta \in \mathbb{R}^3$, but $dim(\{\beta\}) = 2$
$\omega \subset \mathbb{R}^n$, but $dim(\omega) = 2$

The principle of Occam's Razor: *"One should always choose the simplest explanation of a phenomenon, the one that requires the fewest leaps of logic."*

• formulation of hypothesis testing from the view of **comparing models (model spaces)**

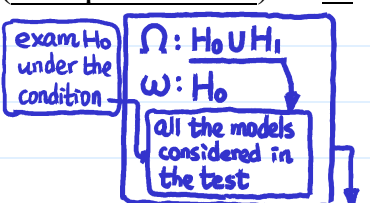➤ a model space ≡ the space spanned by columns of some $X$ (model matrix)

↪ **impose restrictions of $\beta$ on $\Omega$**

➤ consider a large model space, $\Omega$, and a smaller model space, $\omega$, where $\omega \subset \Omega$ (i.e., $\omega$ represents a subset/subspace of $\Omega$). Suppose dimension (# of parameters) of $\Omega$ is $p$ and $dim(\omega) = q$, where $p > q$.

**contains $\Omega$** ← **$c = 0$ or not**

**exam $H_0$ under the condition**

$\Omega: H_0 \cup H_1$
$\omega: H_0$

**all the models considered in the test**

Examples:   **dim = 3**       **dim = 2**       **dim = 2**

$X_\Omega = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ & \cdots & \\ 1 & x_{n1} & x_{n2} \end{bmatrix}$   Q1: $X_\omega = \begin{bmatrix} 1 & x_{12} \\ 1 & x_{22} \\ & \cdots \\ 1 & x_{n2} \end{bmatrix}$   Q2: $X_\omega = \begin{bmatrix} 1 & x_{11} + x_{12} \\ 1 & x_{21} + x_{22} \\ & \cdots \\ 1 & x_{n1} + x_{n2} \end{bmatrix}$

**should also pay attention to $\Omega$. e.g., Is $\Omega$ too simple?**

➤ to answer "which of the model spaces is more adequate" in statistical language $\Rightarrow$ perform the test $H_0$: $\omega$ ($A\beta = c$) v.s. $H_1$: $\Omega \backslash \omega$

**Recall. Example in LNp. 1-15**

- a <u>geometric view</u> of <u>$H_0$</u>: $\omega$   v.s.   <u>$H_1$</u>: $\Omega \backslash \omega$

subspace

$\underline{\Omega \cap \omega^\perp}$

the <u>space</u> that
(1) $\subset \Omega$
(2) $\perp \omega$
(dim = $p-q$)

$\hat{\varepsilon}_\omega$ (($n-q$)-dim)

$Y$ ($n$-dim)

$\hat{\varepsilon}_\Omega$ (($n-p$)-dim)

$\Omega^\perp$ (dim = $n-p$)

$\omega^\perp$

$\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega$:
規律 in $\Omega$
隨机 in $\omega$

∵ 畢氏定理

$$\left\| \hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega \right\|^2$$
$$= \left\| \hat{\varepsilon}_\omega \right\|^2 - \left\| \hat{\varepsilon}_\Omega \right\|^2$$
$$= RSS_\omega - RSS_\Omega$$

$\hat{Y}_\Omega$ ($p$-dim)

$\hat{Y}_\omega$ ($q$-dim)

$0$

(a)

(b) $\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega = \hat{Y}_\Omega - \hat{Y}_\omega$ (($p-q$)-dim)

$\Omega$ (dim = $p$)

$\omega$ (dim = $q$)

cf. Under $H_1$, $Y \sim N(\underline{X_\Omega \beta_\Omega}, \sigma^2 I)$
$Y \sim N(\underline{X_\omega \beta_\omega}, \sigma^2 I)$

- $\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega = \hat{Y}_\Omega - \hat{Y}_\omega$   orthogonal projection matrix
$= (\underline{H_\Omega - H_\omega})Y$

- $(H_\Omega - H_\omega)^2 = H_\Omega - H_\omega$
- $(H_\Omega - H_\omega)^T = H_\Omega - H_\omega$
- $H_\Omega H_\omega = H_\omega H_\Omega = H_\omega$
- $(H_\Omega - H_\omega)(I - H_\Omega)$
$= (I - H_\Omega)(H_\Omega - H_\omega) = 0$
$\Rightarrow \hat{\varepsilon}_\Omega \perp \hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega$

Under <u>$H_0$</u> (<u>null hypothesis</u> $\underline{\omega}$):

By (N1) $\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega \; [= (H_\Omega - H_\omega)Y] \sim N(0, (H_\Omega - H_\omega)\sigma^2)$
<u>not</u> zero vector <u>under</u> $H_1$

By (N6) $RSS_\omega - RSS_\Omega \; [= (\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega)^T(\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega)] \sim \sigma^2 \chi^2_{p-q}$
noncentral chi-square under $H_1$   $(H_\Omega - H_\omega)^- = H_\Omega - H_\omega$

By (N5) $RSS_\omega - RSS_\Omega$ is <u>independent</u> of $RSS_\Omega$
also hold under $H_1$   ∵ $\hat{\varepsilon}_\Omega [= (I - H_\Omega)Y]$ indep. $\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega$

- <u>eigenvalues</u> of $H_\Omega - H_\omega$
are <u>either 0 or 1</u>;
# of 1's = $p-q$;
# of 0's = $n-(p-q)$

- <u>How the geometric view</u> related to <u>test statistic</u> of <u>$H_0$</u>: $\omega$ vs. $H_1$: $\Omega \backslash \omega$?

$\hat{\sigma}^2 = \frac{RSS_\Omega}{n-p}$

unit=? = (unit of $y_i$'s)$^2$

➢ if $RSS_\omega - RSS_\Omega$ is <u>small</u>, $\underline{\omega}$ is a <u>more adequate</u> model relative to $\underline{\Omega}$

➢ suggest $(RSS_\omega - RSS_\Omega)/RSS_\Omega$, where the <u>denominator</u> is used for "scaling" → Q: Why divided by $RSS_\Omega$? Why not divided by $RSS_\omega$? Ans: ① orthogonality ② central $\chi^2$ under $\omega$

➢ Q: What's the <u>scale</u> for $(RSS_\omega - RSS_\Omega)/RSS_\Omega$? ← use <u>null distribution</u> to decide.
i.e., <u>how small</u> is small? <u>how large</u> is large?   $(\sim \frac{p-q}{n-p} F_{p-q, n-p}$ under $H_0)$

subspace

➢ $\underline{\omega}$ can be any of the <u>form</u> $H_0$: $A\beta = 0 \; (\Rightarrow 0 \in \omega)$   $(\omega \ni X\beta = 0$ when $\beta = 0)$
$(p-q) \times 1$   $n \times 1$

➢ <u>generalization</u> to $\underline{\omega}$ of the <u>form</u> $H_0$: $A\beta = c$, where $c \neq 0$, is <u>achievable</u>;
cf.   $(p-q) \times 1$   $p \times 1$

subset

- however, $0 \notin \omega$, and

$\hat{Y}_\omega = \hat{Y}_\omega^* + g_k$
offset

- $\hat{Y}_\omega \perp \hat{\varepsilon}_\omega^*$ does <u>not</u> hold in this case. (but $\hat{Y}_\omega^* \perp \hat{\varepsilon}_\omega^*$)
$n \times 1$

$Y^* = Y - g_k$
$Y^*$ (new $Y$)
$\hat{Y}_\omega^*$
$\hat{\varepsilon}_\omega^*$

$\hat{\varepsilon}_\omega$ (($n-q$)-dim)
$Y$ ($n$-dim)
$\hat{\varepsilon}_\Omega$ (($n-p$)-dim)

$\sqrt{RSS_\Omega}$

parallel
$A\beta = c$
$A\beta = 0$
$g_k$
new $0$
$\theta$

how large is large?

$\hat{Y}_\omega$
$0$
$\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega$ (($p-q$)-dim)

$\Omega$ (dim = $p$)
$\omega$ (dim = $q$)

$\sqrt{RSS_\omega - RSS_\Omega}$

$$\left\| \hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega \right\|^2$$
$$= \left\| \hat{\varepsilon}_\omega \right\|^2 - \left\| \hat{\varepsilon}_\Omega \right\|^2$$
$$= RSS_\omega - RSS_\Omega$$

$$[\tan(\theta)]^2 = \frac{RSS_\omega - RSS_\Omega}{RSS_\Omega}$$

❖ **Reading**: Faraway (2005, 1$^{st}$ ed.), 3.1;     ❖ **Futher reading**: D&S, 21.1, 21.2, 21.3, 21.4