

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

$$\varepsilon \sim \underline{N(0, \sigma^2 I)}$$

➤ model: $Y = X\beta + \varepsilon, \varepsilon \sim \underline{N(0, \sigma^2 I)}$

$$Y \sim \underline{N(X\beta, \sigma^2 I)}$$

- **Q:** what does the model describe? e.g.,

$$y_x = \beta_0 + \beta_1 x + \varepsilon_x, \varepsilon_x \text{'s} \sim \text{i.i.d. } \underline{N(0, \sigma^2)}$$

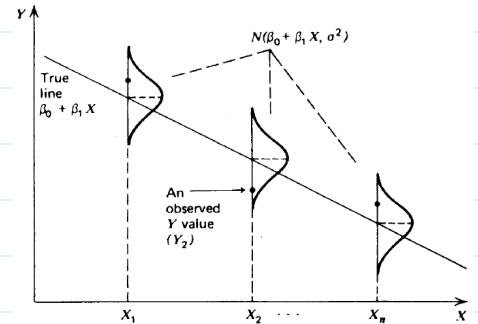
$$\Rightarrow E(y_x) = \beta_0 + \beta_1 x$$

$$\Rightarrow y_x \text{'s are independent and } y_x \sim \underline{N(\beta_0 + \beta_1 x, \sigma^2)} \text{ at } x = x_i, i = 1, \dots, n.$$

- **Q:** how should the data generated from the model look like?
- **Q:** when would it be appropriate to impose the inference based on this regression model on the underlying true model? Can we use it when these exist clear differences between the two models, e.g., what if y is a discrete quantitative measurement?

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

George Box: "all models are wrong, but some are useful"



NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

➤ **Q:** why Normal?

- CLT \Rightarrow when random error is a sum of many small random disturbances
- bell shape curve is common
- from the viewpoint of approximation
- good mathematical/statistical properties

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

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

➤ **Q:** under what conditions, the Normality assumption is inappropriate?

- qualitative response
- quantitative discrete response
with only few possible outcomes
- skewed error
- heavy tail error

* Some properties of (multivariate) Normal distribution

(N1). linear transformation of Normal is still Normal

$$\underline{Z} \sim N(\underline{\mu}, \underline{\Sigma}) \Rightarrow \underline{AZ+c} \sim N(\underline{A}\underline{\mu}+\underline{c}, \underline{A}\underline{\Sigma}\underline{A}^T)$$

(N2). when 1st and 2nd moments are given, the Normal distribution is specified

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

(N4). $\underline{Z} \sim N(\underline{\mu}, \underline{\Sigma}), W_1=\underline{A}_1\underline{Z}, W_2=\underline{A}_2\underline{Z} \Rightarrow$ W_1, W_2 are independent iff $\underline{A}_1\underline{\Sigma}\underline{A}_2^T=\underline{0}$

(N5). $\underline{Z} \sim N(\underline{\mu}, \underline{\Sigma}), W_1=\underline{A}_1\underline{Z}, W_2=\underline{A}_2\underline{Z}, \dots, W_k=\underline{A}_k\underline{Z}$, and $\text{cov}(W_i, W_j)=\underline{0}$ for $i \neq j$
 \Rightarrow $W_1^T W_1, W_2^T W_2, \dots, W_k^T W_k$ are mutually independent

(N6). \underline{Z} : an $n \times 1$ random vector and $\underline{Z} \sim N(\underline{\mu}, \underline{\Sigma})$, then

- $(\underline{Z}-\underline{\mu})^T \underline{\Sigma}^{-1}(\underline{Z}-\underline{\mu}) \sim \chi_{n-2}^2$ if $\underline{\Sigma}$ is non-singular
- $(\underline{Z}-\underline{\mu})^T \underline{\Sigma}^-(\underline{Z}-\underline{\mu}) \sim \chi_{r-2}^2$ if $\underline{\Sigma}$ is singular and has rank $r (< n)$,
 where $\underline{\Sigma}^-$ is a generalized inverse of $\underline{\Sigma}$, (i.e., $\underline{\Sigma}\underline{\Sigma}^-\underline{\Sigma}=\underline{\Sigma}$)

• Some properties of linear models when $\underline{\epsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$:

- distribution of $\underline{Y} [= \underline{X}\underline{\beta} + \underline{\epsilon}] \sim N(\underline{X}\underline{\beta}, \sigma^2 \underline{I})$
- distribution of $\hat{\underline{\beta}} [= (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}] \sim N(\underline{\beta}, (\underline{X}^T \underline{X})^{-1} \sigma^2)$

NTHU STAT 5410, 2022, Lecture Notes
 made by S.-W. Cheng (NTHU, Taiwan)

➤ distribution of $\hat{\underline{\epsilon}} [= (\underline{I}-\underline{H})\underline{Y}=(\underline{I}-\underline{H})\underline{\epsilon}] \sim N(\underline{0}, (\underline{I}-\underline{H})\sigma^2)$, which has a singular covariance matrix $\underline{I}-\underline{H}$ with rank $n-p$ (Note: $\text{dim}(\hat{\underline{\epsilon}})=n-p$)

➤ distribution of $RSS [= (n-p)\hat{\sigma}^2 = \hat{\underline{\epsilon}}^T \hat{\underline{\epsilon}} = \underline{\epsilon}^T (\underline{I}-\underline{H})\underline{\epsilon}] \sim \sigma^2 \chi_{n-p}^2$

➤ distribution of $\hat{\underline{Y}} [= \underline{X}\hat{\underline{\beta}} = \underline{H}\underline{Y}] \sim N(\underline{X}\underline{\beta}, \underline{H}\sigma^2)$, which has a singular covariance matrix with rank p (Note: $\text{dim}(\hat{\underline{Y}})=p$)

➤ $\hat{\underline{\beta}}$ is independent of $\hat{\sigma}^2$ (Note: $\text{cov}((\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}, (\underline{I}-\underline{H})\underline{Y})=\underline{0}$)

➤ $\hat{\underline{Y}}$ is independent of $\hat{\underline{\epsilon}}$ (Note: $\text{cov}(\underline{H}\underline{Y}, (\underline{I}-\underline{H})\underline{Y})=\underline{0}$)

➤ distribution of prediction for a new set of predictors, $\underline{x}_0 = (g_1(x_{10}, \dots, x_{m0}), \dots, g_p(x_{10}, \dots, x_{m0}))^T$
 model: $y = \sum_{j=1}^p \beta_j \cdot g_j(x_1, \dots, x_m) + \epsilon$

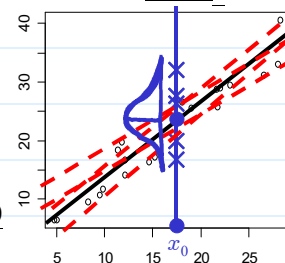
- mean response v.s. future observation (Q: what different?)
 - Example: average yield when $\underline{x}=\underline{x}_0$? and tomorrow's yield when $\underline{x}=\underline{x}_0$?
 - same predicted value $\underline{x}_0^T \hat{\underline{\beta}}$, but different distributions

▪ distribution of prediction error for mean response at \underline{x}_0

$$\underline{x}_0^T \hat{\underline{\beta}} - \underline{x}_0^T \underline{\beta} \sim N(\underline{0}, (\underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0) \sigma^2)$$

▪ distribution of prediction error for future observations at \underline{x}_0

$$\underline{x}_0^T \hat{\underline{\beta}} - (\underline{x}_0^T \underline{\beta} + \epsilon) \sim N(\underline{0}, (\underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0 + \underline{1}) \sigma^2)$$



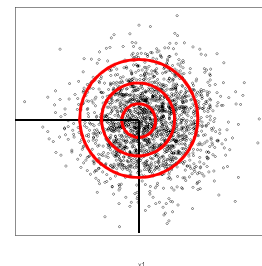
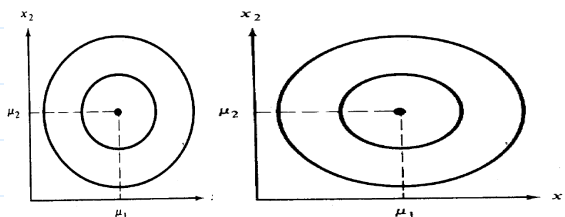
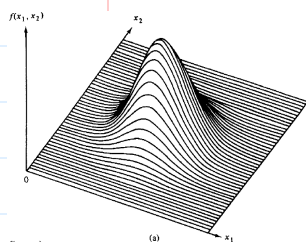
➤ Example: bivariate normal $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}\right)$

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

joint pdf

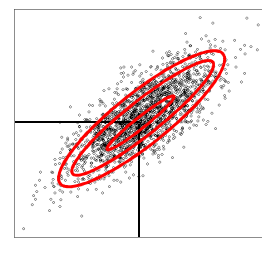
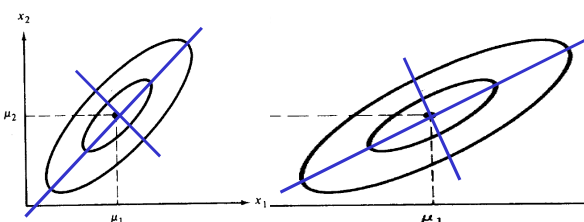
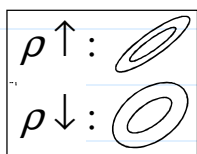
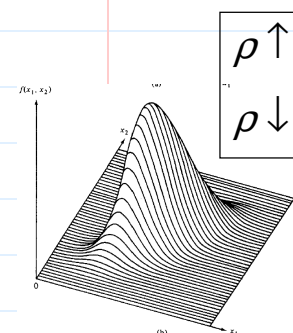
contour lines of the pdf

data generated from the pdf



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

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



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 $(x-\mu)^T \Sigma^{-1} (x-\mu) = c$

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

hypothesis testings (for β)

• **Q:** What questions is a hypothesis testing about β trying to answer?

examples: full model: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$

Q1: $\beta_1 = 0?$
 $\Rightarrow y = \beta_0 + \beta_2 x_2 + \epsilon$

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

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

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

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 \equiv the space spanned by columns of some X (model matrix)

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

Examples:

$$X_{\Omega} = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ \dots & \dots & \dots \\ 1 & x_{n1} & x_{n2} \end{bmatrix}$$

Q1:

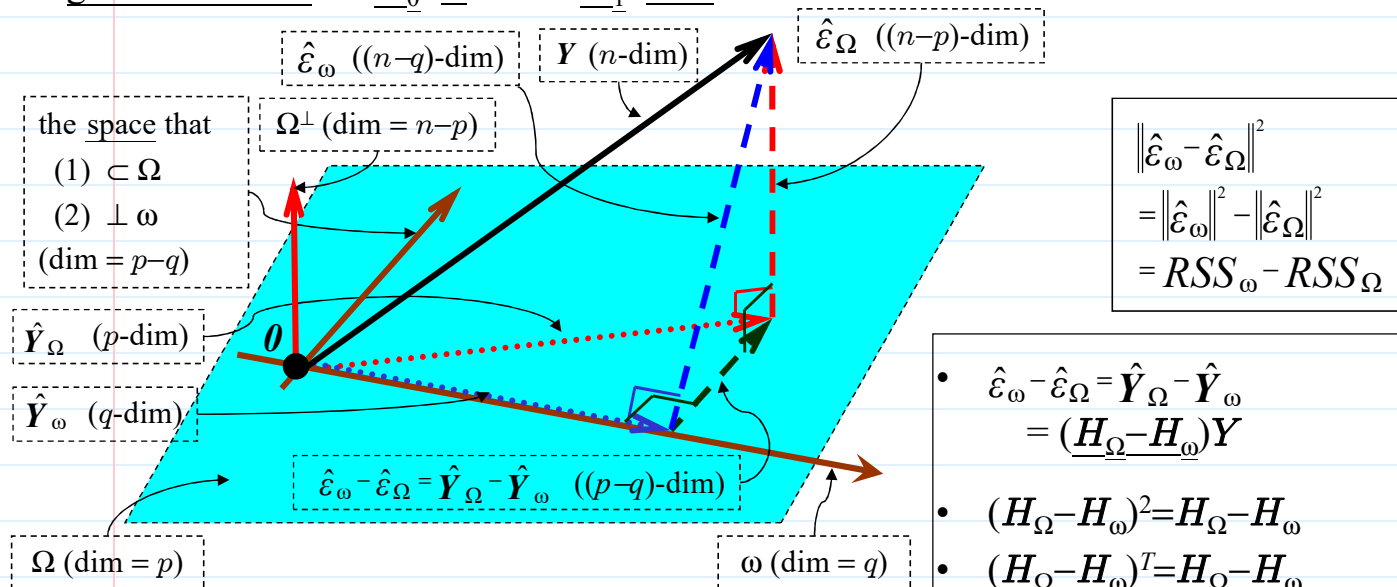
$$X_{\omega} = \begin{bmatrix} 1 & x_{12} \\ 1 & x_{22} \\ \dots & \dots \\ 1 & x_{n2} \end{bmatrix}$$

Q2:

$$X_{\omega} = \begin{bmatrix} 1 & x_{11} + x_{12} \\ 1 & x_{21} + x_{22} \\ \dots & \dots \\ 1 & x_{n1} + x_{n2} \end{bmatrix}$$

➤ 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. $H_1: \Omega \setminus \omega$

• a geometric view of $H_0: \omega$ v.s. $H_1: \Omega \setminus \omega$



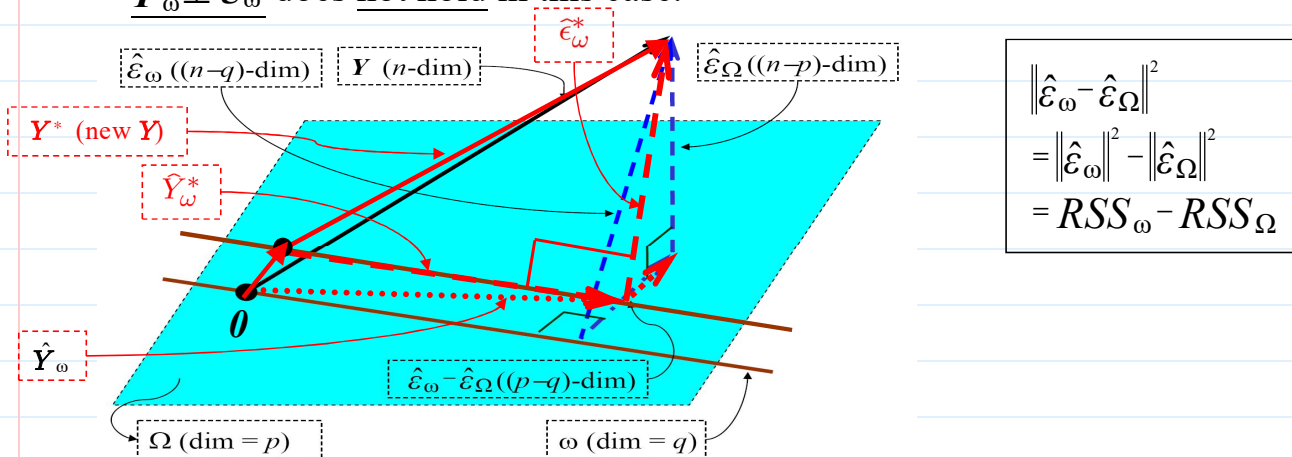
Under H_0 (null hypothesis ω):

- $\hat{e}_\omega - \hat{e}_\Omega [= (H_\Omega - H_\omega)Y] \sim N(\mathbf{0}, (H_\Omega - H_\omega)\sigma^2)$
- $RSS_\omega - RSS_\Omega [= (\hat{e}_\omega - \hat{e}_\Omega)^T(\hat{e}_\omega - \hat{e}_\Omega)] \sim \sigma^2 \chi^2_{p-q}$
- $RSS_\omega - RSS_\Omega$ is independent of RSS_Ω

- $\hat{e}_\omega - \hat{e}_\Omega = \hat{Y}_\Omega - \hat{Y}_\omega = (H_\Omega - H_\omega)Y$
- $(H_\Omega - H_\omega)^2 = H_\Omega - H_\omega$
- $(H_\Omega - H_\omega)^T = H_\Omega - H_\omega$
- $H_\Omega H_\omega = H_\omega H_\Omega = H_\omega$
- $(H_\Omega - H_\omega)(I - H_\Omega) = (I - H_\Omega)(H_\Omega - H_\omega) = \mathbf{0}$
 $\Rightarrow \hat{e}_\Omega \perp \hat{e}_\omega - \hat{e}_\Omega$
- eigenvalues of $H_\Omega - H_\omega$ are either 0 or 1;
 # of 1's = $p-q$;
 # of 0's = $n-(p-q)$

NTHU STAT 5410, 2022, Lecture Notes
 made by S.-W. Cheng (NTHU, Taiwan)

- How the geometric view related to test statistic of $H_0: \omega$ vs. $H_1: \Omega \setminus \omega$?
 - if $RSS_\omega - RSS_\Omega$ is small, ω is a more adequate model relative to Ω
 - suggest $(RSS_\omega - RSS_\Omega) / RSS_\Omega$, where the denominator is used for "scaling"
 - Q: What's the scale for $(RSS_\omega - RSS_\Omega) / RSS_\Omega$? i.e., how small is small? how large is large?
 - ω can be any of the form $H_0: A\beta = \mathbf{0}$ ($\Rightarrow \mathbf{0} \in \omega$)
 - generalization to ω of the form $H_0: A\beta = \mathbf{c}$, where $\mathbf{c} \neq \mathbf{0}$, is achievable;
 - however, $\mathbf{0} \notin \omega$, and
 - $\hat{Y}_\omega \perp \hat{e}_\omega$ does not hold in this case.



• Example 1: test of all predictors

➤ **Q:** are any of the predictors g_i 's useful in predicting the response?

▪ $\Omega: y = \beta_0 + \beta_1 g_1 + \dots + \beta_{p-1} g_{p-1} + \epsilon$, $\dim(\Omega)=$, $df_{\Omega}=$

▪ $\omega:$, $\dim(\omega)=$, $df_{\omega}=$

▪ $H_0:$ $H_1:$

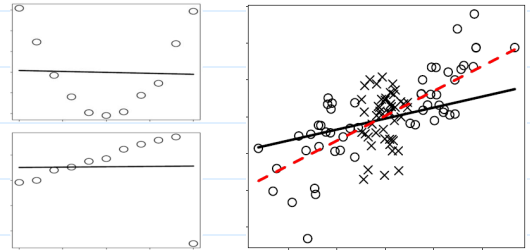
▪ $RSS_{\Omega}:$ $RSS_{\omega}:$

▪ (the overall F) $F = \frac{(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})}{RSS_{\Omega} / df_{\Omega}}$

➤ **Q:** What's the "meaning" of H_0 ? Let's consider the following two questions:

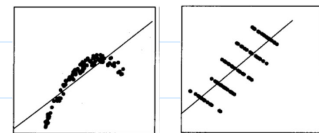
▪ If H_0 is not rejected, what can you conclude? is it the end of the analysis?

Ans: No. Check assumptions, such as linearity, outlier, or if enough data are collected, Do not conclude too soon that no real relationship exist between Y and X_1, \dots, X_p .



▪ If H_0 is rejected, does it mean the alternative model is the best choice?

Ans: No. Check if some predictors can be dropped, if other predictors might be added, ...



• Example 2: testing just one predictor

➤ **Q:** Can one particular predictor, say $g_i(x)$, be dropped from the model?

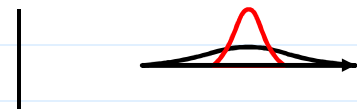
▪ $\Omega: y = \beta_0 + \dots + \beta_i g_i + \dots + \beta_{p-1} g_{p-1} + \epsilon$, $\dim(\Omega)=$, $df_{\Omega}=$

▪ $\omega: y = \beta_0 + \dots + \beta_i g_i + \dots + \beta_{p-1} g_{p-1} + \epsilon$, $\dim(\omega)=$, $df_{\omega}=$

▪ $H_0:$ $H_1:$

▪ $F = [(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})] / (RSS_{\Omega} / df_{\Omega}) \sim F_{df_{\omega} - df_{\Omega}, df_{\Omega}}$

➤ alternative method t-test: $t_i = \hat{\beta}_i / se(\hat{\beta}_i) \sim t_{n-p}$ [Note. $t_i^2 \sim F_{1, n-p}$, and $t_i^2 = F$]



➤ **Q:** What is the "meaning" of H_0 ? It seems only β_i appears in null, does H_0 say anything about other β_j 's, where $j \neq i$?

Note. all g_j 's, where $j \neq i$, are included in both ω and Ω .

Ans: when all other predictors are included in the model, whether g_i is helpful in interpreting the response variation.

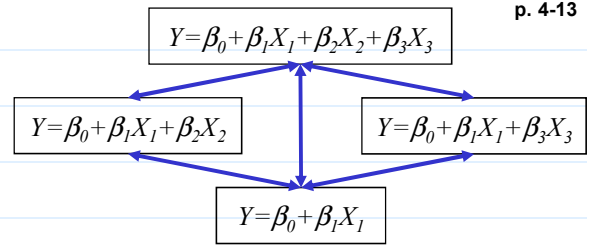
➤ **Q:** When "other predictors" are changed, can we always get the same result for the test of g_i ? **Ans. NO**, but why?

➤ **Q:** When can rejecting/accepting $H_0: \beta_i = 0$ "almost irrelevant" to whether other predictors appear in the models or not?

Hint. what will happen if g_i is orthogonal to all g_j 's, where $j \neq i$? under this condition, $\hat{\beta}_i$ independent of all $\hat{\beta}_j$'s? try give it a geometric interpretation.

• Example 3: testing a pair of predictors

- **Q:** Suppose the t -tests for $\underline{\beta}_j$ and $\underline{\beta}_k$ are both insignificant, can you remove both g_j and g_k from the model? when can and when cannot? and why? (**Hint:** what's the null in the 2 t -tests?)



- **Q:** What combinations of acceptance/rejection you will see in these tests?

- **Q:** Can two particular predictors, say g_j and g_k , be dropped from the model?

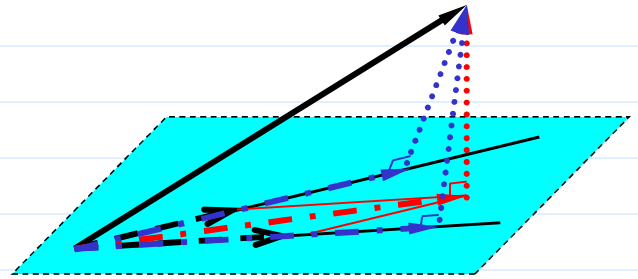
- $\Omega: y = \beta_0 + \dots + \beta_j g_j + \dots + \beta_k g_k + \dots + \epsilon$, $\dim(\Omega) =$, $df_{\Omega} =$

- $\omega: y = \beta_0 + \dots + \beta_j g_j + \dots + \beta_k g_k + \dots + \epsilon$, $\dim(\omega) =$, $df_{\omega} =$

- $H_0:$ $H_1:$

- $F = [(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})] / (RSS_{\Omega} / df_{\Omega}) \sim F_{df_{\omega} - df_{\Omega}, df_{\Omega}}$

- **Q:** When the data accept $H_{0,i}: \beta_j = 0$ and $H_{0,k}: \beta_k = 0$, but reject $H_0: \beta_j = \beta_k = 0$, how can you explain the contradictory results? how is it related to orthogonality and collinearity?



- It can be generalized to more than two predictors. How? (**exercise**)

• Example 4: testing a subspace/subset ω

- **Q:** how to test $H_0: \beta_j + \beta_k = 1$?

- $\Omega: y = \beta_0 + \dots + \beta_j g_j + \dots + \beta_k g_k + \dots + \epsilon$, $\dim(\Omega) =$, $df_{\Omega} =$

- $\omega:$, $\dim(\omega) =$, $df_{\omega} =$

offset

- $F = [(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})] / (RSS_{\Omega} / df_{\Omega}) \sim F_{df_{\omega} - df_{\Omega}, df_{\Omega}}$

- **Q:** how to test $H_0: \beta_j = \beta_k$?

- $\Omega: y = \beta_0 + \dots + \beta_j g_j + \dots + \beta_k g_k + \dots + \epsilon$, $\dim(\Omega) =$, $df_{\Omega} =$

- $\omega:$, $\dim(\omega) =$, $df_{\omega} =$

- $F = [(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})] / (RSS_{\Omega} / df_{\Omega}) \sim F_{df_{\omega} - df_{\Omega}, df_{\Omega}}$

- **Q:** how to test $H_0: \beta_j = c$, c : a known constant, say $\beta_j = 10$?

- $\Omega: y = \beta_0 + \dots + \beta_j g_j + \dots + \beta_{p-1} g_{p-1} + \dots + \epsilon$, $\dim(\Omega) =$, $df_{\Omega} =$

- $\omega:$, $\dim(\omega) =$, $df_{\omega} =$

- $F = [(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})] / (RSS_{\Omega} / df_{\Omega}) \sim F_{df_{\omega} - df_{\Omega}, df_{\Omega}}$

- alternative method t -test: $t_j = (\hat{\beta}_j - c) / \text{se}(\hat{\beta}_j) \sim t_{n-p}$

- **Q:** Can we apply the method to test $H_0: \beta_j \beta_k = 1$?

- $\Omega: y = \beta_0 + \dots + \beta_j g_j + \dots + \beta_k g_k + \dots + \epsilon$

- $\omega: y = \beta_0 + \dots + \beta_j g_j + \dots + (1/\beta_j) g_k + \dots + \epsilon$

- Some note & concerns about hypothesis testing

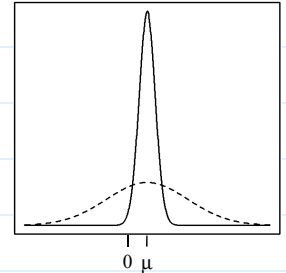
- The previous testing method can be applied to $H_0: \underline{A}\underline{\beta}=\underline{c}$, where \underline{A} is a known $(p-q)\times p$ matrix of rank $p-q$, and \underline{c} is a known $(p-q)\times 1$ vector.
examples:

Q: what are ω and Ω ?

- **Q:** Suppose (1) the model is correct and (2) the estimators of $\underline{\beta}$ are mutually independent. When $H_0: \underline{\beta}_i=0$ is accepted, does it really mean that $\underline{\beta}_i$ is exactly zero?

e.g.: $y_i = \mu + \epsilon_i$; $\epsilon_i \sim$ i.i.d. $N(0, \sigma^2)$; $\mu \approx 0$, but not zero

$$H_0 : \mu = 0$$



Note: that's why we usually don't say "accept H_0 ", but say "sample size isn't large enough to reject H_0 ".

- When sample size, n , is much larger than the number of parameters, p , it's very possible that every tests are significant (even though R^2 is very low).

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- Statistical significance may not be equivalent to practical/physical significance. p. 4-16
example:

- **Q:** why inequivalent? **Hint:** what are the numerator & denominator in the t-test? Does the denominator represents a scale of physical significance?)
- for datasets with large n , it is easy to get statistically significant results on $\underline{\beta}_i$'s, but the magnitudes of some (all) $\underline{\beta}_i$'s may be quite small and therefore, not physically important.

- The inference depends on the correctness of the model $\Omega : Y = X\underline{\beta} + \epsilon$ we use. The assumptions about the model can be checked, but there will be always some element of doubt. (**Q:** what you can do?)

- The data may suggest more than one possible models which may lead to contradictory results, e.g, when strong collinearity exists. (**Q:** what you can do?)

- What is the true significant level of several tests, each with significant level α ?

- ❖ **Reading:** Faraway (2005, 1st ed.), 3.2
- ❖ **Further reading:** D&S, 9.1