

Analysis of Variance

overall F-test
 $(\omega) H_0: y = \beta_0 + \epsilon$
 $(\Omega) H_A: y = \beta_0 + \sum_{j=1}^k \beta_j x_j + \epsilon$

- The total variation in y , i.e., corrected total sum of squares, $\|y - \bar{y}\mathbf{1}\|^2$

$CTSS = \sum_{i=1}^N (y_i - \bar{y})^2 = \mathbf{y}^T \mathbf{y} - N\bar{y}^2$, can be decomposed into two parts

(Analysis of Variance (ANOVA))

規律 (regularity)
 隨機 (stochastic)

RSS under Ω

RSS under ω

$CTSS = RegrSS + RSS$

$-\hat{y}_i + \hat{y}_i$

where $RSS = \text{Residual sum of squares} = \sum (y_i - \hat{y}_i)^2 = (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta})$,

$RegrSS = \text{Regression sum of squares} = \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 = \hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2$.

mean of \hat{y}_i 's = \bar{y}

ANOVA Table

y_i 's $\rightarrow y$, \hat{y}_i 's $\rightarrow \mathbf{X}\hat{\beta}$, $\bar{y} \rightarrow \bar{y}\mathbf{1}$

check the graph in LNp.2-20

Source of variation

| Source | Degrees of Freedom | Sum of Squares | Mean Squares |
|------------|--|---|---|
| regression | k $\leftarrow \beta_1, \dots, \beta_k$ | $\hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2$ | $(\hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2) / k$ |
| residual | $N - (k + 1)$ | $(\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta})$ | $(\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta}) / (N - k - 1)$ |
| total | $N - 1$ | $\mathbf{y}^T \mathbf{y} - N\bar{y}^2$ | |

$F\text{-statistic} = \frac{\text{Mean Square Regression}}{\text{Mean Square Residual}}$

Explanatory Power of the Model

規律 (regularity)

隨機 (stochastic)

- $R^2 = \frac{RegrSS}{CTSS} = 1 - \frac{RSS}{CTSS}$ measures of the proportion of variation in y explained by the fitted model. R is called the multiple correlation coefficient.

$\uparrow \text{cor}(y, \hat{y}) \leftarrow \text{check the graph in LNp.2-20}$

- Adjusted R^2 :

$R_a^2 = 1 - \frac{RSS}{\frac{CTSS}{N-1}} = 1 - \left(\frac{N-1}{N-k-1} \right) \frac{RSS}{CTSS} \leq R^2$

$\uparrow \geq 1$

\therefore RSS always decreases

- When an additional predictor is included in the regression model, R^2 always increases. This is not a desirable property for model selection. However, R_a^2 may decrease if the included variable is not an informative predictor.

Usually R_a^2 is a better measure for comparing different model fits.

Q: When can R^2 reach 1?
 Ans. $\hat{y}_i = y_i$, only when data have no replicates & # of parameters $\beta = \#$ of obs

total df available N \leftarrow sample size
 \rightarrow df available for studying β
 \rightarrow df available for estimating σ^2

Want to have a final fitted model such that
 1. # of parameters: small
 2. RSS: small

Testing significance of coefficients : t-Statistic

Examine whether $\beta_j=0$ when other effects β_i 's (g_i 's), $i \neq j$, are still in the model

- To test the null hypothesis $H_0 : \beta_j = 0$ against the alternative hypothesis $H_A : \beta_j \neq 0$ under the full model, use the test statistic

Note: collinearity

$$t_j^2 = F$$

$$t_j = \frac{\hat{\beta}_j - 0}{s.d.(\hat{\beta}_j)}$$

↖ null model

$$(\omega) H_0 : y = \beta_0 + \dots + \beta_{j-1}g_{j-1} + \beta_{j+1}g_{j+1} + \dots + \epsilon$$

↗ alternative model

$$(\Omega) H_A : y = \beta_0 + \sum_{i=1}^k \beta_i g_i + \epsilon$$

- The higher the value of $|t_j|$, the more significant is the coefficient.
- In practice, if p -value is less than $\alpha = 0.05$ or 0.01 , H_0 is rejected.
- Confidence Interval** : $100(1 - \alpha)\%$ confidence interval for β_j is given by

$$\hat{\beta}_j \pm t_{N-(k+1), \frac{\alpha}{2}} \times s.d.(\hat{\beta}_j)$$

estimate critical value s.d.(estimate)

$H_0 : \beta_j = \beta_j^*$

$$\left| \frac{\hat{\beta}_j - \beta_j^*}{s.d.(\hat{\beta}_j)} \right| \leq t_{N-(k+1), \frac{\alpha}{2}}$$

acceptance region

where $t_{N-k-1, \frac{\alpha}{2}}$ is the upper $\alpha/2$ point of the t distribution with $N - k - 1$ degrees of freedom.

df_{RSS}

If the confidence interval for β_j does not contain 0, then H_0 is rejected.

Analysis of Air Pollution Data

check graphs in LNp.2-18

| Predictor | Coef | SE Coef | T | P |
|----------------------|------------|------------|-------|---------|
| ① β_0 Constant | 1332.7 | 291.7 | 4.57 | 0.000 ✓ |
| ② JanTemp | -2.3052 | 0.8795 | -2.62 | 0.012 ✓ |
| ③ JulyTemp | -1.657 | 2.051 | -0.81 | 0.424 |
| ④ RelHum | 0.407 | 1.070 | 0.38 | 0.706 |
| ⑤ Rain | 1.4436 | 0.5847 | 2.47 | 0.018 ✓ |
| ⑥ Educatio | -9.458 | 9.080 | -1.04 | 0.303 |
| ⑦ PopDensi | 0.004509 | 0.004311 | 1.05 | 0.301 |
| ⑧ %NonWhit | 5.194 | 1.005 | 5.17 | 0.000 ✓ |
| ⑨ %WC | -1.852 | 1.210 | -1.53 | 0.133 |
| ⑩ pop | 0.00000109 | 0.00000401 | 0.27 | 0.788 |
| ⑪ pop/hous | -45.95 | 39.78 | -1.16 | 0.254 |
| ⑫ income | -0.000549 | 0.001309 | -0.42 | 0.677 |
| ⑬ logHC | -53.47 | 35.39 | -1.51 | 0.138 |
| ⑭ logNOx | 80.22 | 32.66 | 2.46 | 0.018 ✓ |
| logSO2 | -6.91 | 16.72 | -0.41 | 0.681 |

Q: Can we remove all insignificant effects simultaneously to simplify the model?
 Ans. No, in general, but OK if orthogonality exists.

how to interpret it?
 It becomes insignificant when other effects are in the model. (cf. graph in LNp.2-18)
 Possible reason: collinearity

S = 34.58 R-Sq = 76.7% R-Sq(adj) = 69.3%

Analysis of Variance

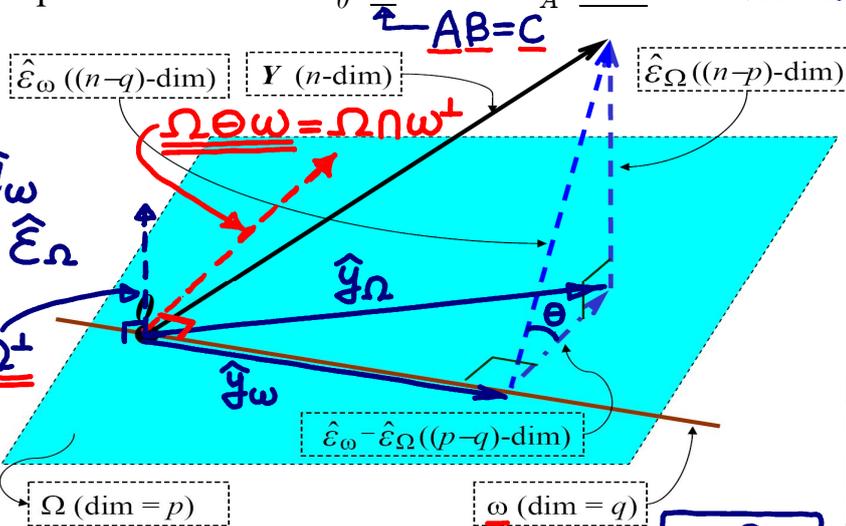
$$R^2 = 1 - \frac{1}{1 + \frac{k}{N-(k+1)} F}$$

| Source | DF | SS | MS | F | P |
|----------------|----|--------|-------|-------|-------|
| Regression | 14 | 173383 | 12384 | 10.36 | 0.000 |
| Residual Error | 44 | 52610 | 1196 | | |
| Total | 58 | 225993 | | | |

59 obs.

- formulation of hypothesis testing from the view of **comparing models** LM, LNp.4-6~16
- a model space \equiv the space spanned by columns of some X **effects**
- consider a large model space, Ω , and a smaller model space, ω , where $\omega \subset \Omega$, i.e., ω represents a subset/a subspace of Ω . Suppose dimension (# of parameters) of Ω is p and $\dim(\omega) = q$, where $p > q$. $\rightarrow df_{\Omega} = n - p, df_{\omega} = n - q$
- to answer “which of the model spaces is more adequate” in statistical language \rightarrow perform the test $H_0: \omega$ v.s. $H_A: \Omega \setminus \omega$ \leftarrow check examples in LNp.2-21&23

$H_0: \beta_j = \beta_j^* = 0$



- θ large $\Rightarrow \hat{E}_{\omega} \approx \hat{E}_{\Omega}$
 $\Rightarrow \hat{y}_{\Omega} \approx \hat{y}_{\omega}$
 \Rightarrow prefer ω
- θ small
 $\Rightarrow \hat{E}_{\omega}$ quite different from \hat{E}_{Ω}
 $\Rightarrow \hat{y}_{\Omega}$ quite different from \hat{y}_{ω}
 \Rightarrow prefer Ω

$$F = \frac{\frac{RSS_{\omega} - RSS_{\Omega}}{p - q}}{\frac{RSS_{\Omega}}{n - p}} \sim F_{p-q, n-p} \text{ (under } \omega \text{)}$$

\leftarrow a special case: t-test

(sequential) ANOVA \leftrightarrow t-tests (LNp.2-23) & F-test (LNp.2-25)

LM, LNp.8-18~19 \square anova($y \sim I + A + B + A:B$), A: 3 levels, B: 4 levels $\Rightarrow A:B: 2 \times 3 = 6$

2 dummy var. 3 dummy var. dummy var. drop one from the full model

- 1) test ω : model 1 ($y \sim I$) against Ω : model 2 ($y \sim I + A$) [$df_{\omega} - df_{\Omega} = 2$]
- 2) test ω : model 2 ($y \sim I + A$) against Ω : model 4 ($y \sim I + A + B$) [$df_{\omega} - df_{\Omega} = 3$]
- 3) test ω : model 4 ($y \sim I + A + B$) against Ω : model 5 ($y \sim I + A + B + A:B$) [$df_{\omega} - df_{\Omega} = 6$]

$$F = \frac{(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})}{RSS_{\text{model 5}} / df_{\text{model 5}}} \sim F_{df_{\omega} - df_{\Omega}, df_{\text{model 5}}}$$

$\hat{\sigma}_{\text{Full model}}^2$ \leftarrow full model \leftarrow cf.

\square invariant to the choice of dummy variables since they generate same ω and Ω

\square ANOVA could have different results when the order of effect sequence is changed, e.g., anova($y \sim I + B + A + A:B$):

$$RSS_{\omega} - RSS_{\Omega} = \|P_{\Omega \setminus \omega} Y\|^2$$

\leftarrow projection

- α) test ω : model 1 ($y \sim I$) against Ω : model 3 ($y \sim I + B$) [$df_{\omega} - df_{\Omega} = 3$]
- β) test ω : model 3 ($y \sim I + B$) against Ω : model 4 ($y \sim I + B + A$) [$df_{\omega} - df_{\Omega} = 2$]
- χ) test ω : model 4 ($y \sim I + B + A$) against Ω : model 5 ($y \sim I + B + A + A:B$) [$df_{\omega} - df_{\Omega} = 6$]

\square anova($y \sim I + A + B + A:B$) and anova($y \sim I + B + A + A:B$) will have **identical** results when orthogonality exists between the three groups of effects: $\text{span}\{d_j^A\}$, $\text{span}\{d_j^B\}$, $\text{span}\{d_{jj}^{A:B}\}$, because in the case, $RSS_{\omega} = RSS_{\Omega}$ would equal for 1) and β), 2) and α), 3) and χ) \rightarrow furthermore, also identical to the drop-one t-tests & F-tests

consider the full model: \leftrightarrow submodels

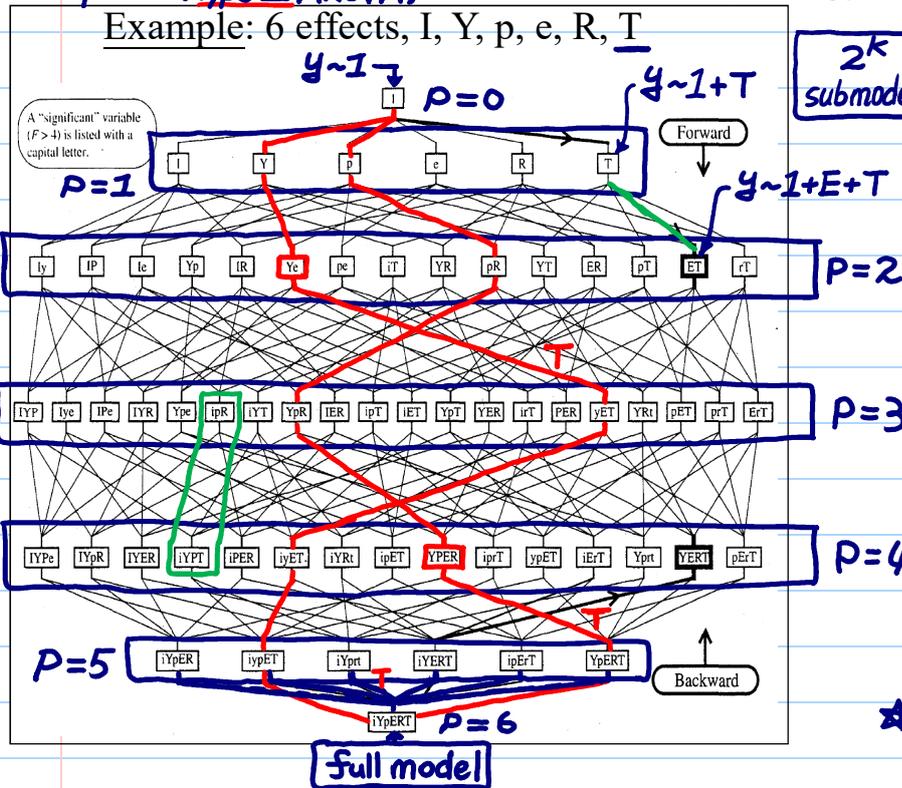
$$y = \beta_0 + \beta_1 g_1(x_1, \dots, x_m) + \beta_2 g_2(x_1, \dots, x_m) + \dots + \beta_k g_k(x_1, \dots, x_m) + \epsilon$$

For $1 \leq i \leq k$, should the term $\beta_i g_i$ be included in the final fitted model?

— : sequential (Type I ANOVA)

— : drop one (Type III ANOVA)

main purpose of performing t- & F-tests



(sub-)model: a model with a subset of all k terms, e.g.,

$$\{1, g_1, g_2\} \rightarrow y \sim 1 + g_1 + g_2$$

$$\{1, g_2, g_4, g_5, g_k\}, \dots \rightarrow y \sim 1 + g_2 + g_4 + g_5 + g_k$$

hierarchical structure of all sub-models (see graph)

- $p = \#$ of terms in a sub-model
- $\#$ of different sub-models = 2^k
- connecting line: model nesting

★ can be generated to "groups" of effects

LM.LNp.5-8~9 → **Orthogonality**

Q: consider the two models:

$$Y = X_1 \beta_1 + \epsilon \quad Y = X_1 \beta_1 + X_2 \beta_2 + \epsilon$$

model 1: $y = \beta_0 + \beta_1 x_1 + \epsilon$, model 2: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$

In general, $\hat{\beta}_1$ in the 2 models are not identical (of course, test $H_0: \beta_1 = 0$ not identical neither) an exception: when x_1 and x_2 are orthogonal

$X_1 (X_1^T X_1)^{-1} X_1^T$ (hat matrix, a projection matrix)

$H X_2 \beta_2 (I - H) X_2 \beta_2$

fitted model = model 1, true model = model 2

What if $X_1^T X_2 = 0$?

$$E(\hat{\beta}_1) = \beta_1 + (X_1^T X_1)^{-1} X_1^T X_2 \beta_2$$

• $Y = X\beta + \epsilon = X_1 \beta_1 + X_2 \beta_2 + \epsilon$, where $\beta = [\beta_1 \ \beta_2]^T$ and $X = [X_1 \ X_2]$ with the property

$X_1^T X_2 = 0 \Rightarrow X_1$ and X_2 are orthogonal $\hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = (X^T X)^{-1} X^T Y = \begin{bmatrix} (X_1^T X_1)^{-1} X_1^T Y \\ (X_2^T X_2)^{-1} X_2^T Y \end{bmatrix}$

$\leftarrow \text{span}\{X_1\} \perp \text{span}\{X_2\}$

$$X^T X = \begin{pmatrix} X_1^T X_1 & X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 \end{pmatrix} = \begin{pmatrix} X_1^T X_1 & 0 \\ 0 & X_2^T X_2 \end{pmatrix} \Rightarrow (X^T X)^{-1} = \begin{pmatrix} (X_1^T X_1)^{-1} & 0 \\ 0 & (X_2^T X_2)^{-1} \end{pmatrix}$$

• Estimation: $\hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T Y$, $\hat{\beta}_2 = (X_2^T X_2)^{-1} X_2^T Y$, and $\hat{\beta}_1, \hat{\beta}_2$ independent \Rightarrow note that $\hat{\beta}_1$ will be the same regardless of whether X_2 is in the model or not (and vice versa). Under model $Y = X_1 \beta_1 + \epsilon \Rightarrow \hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T Y$

Q: what if only two predictors, say some x_i in X_1 and some x_j in X_2 , are orthogonal?

• Randomization: In an exp't, suppose that true model is $Y = X\beta + Z\gamma + \epsilon$, but Z cannot be measured or may not even be suspected $\Rightarrow E(\hat{\beta}) = \beta + (X^T X)^{-1} X^T Z \gamma \Rightarrow$

Q: what's the best way of controlling X to make X and Z as orthogonal as possible?

• Generalization.

$$Y = \beta_0 \mathbf{1} + X_1 \beta_1 + X_2 \beta_2 + \dots + X_k \beta_k + \epsilon$$

$$W_0 = \text{span}\{\mathbf{1}\} \quad W_1 = \text{span}\{X_1\} \quad W_2 = \text{span}\{X_2\} \quad W_k = \text{span}\{X_k\}$$

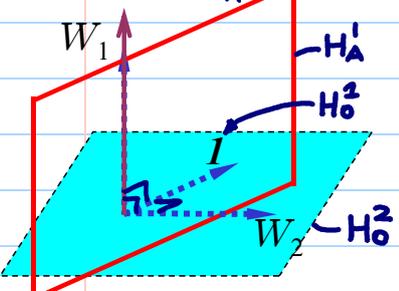
Suppose that $W_0 \perp W_1 \perp W_2 \perp \dots \perp W_k \Rightarrow X^T X =$

$$\frac{RSS_{H_0^1} - RSS_{H_A^1}}{RSS_{H_0^2} - RSS_{H_A^2}} = \frac{\|P_{W_1}(y)\|^2}{\|P_{W_2}(y)\|^2}$$

orthogonal

$$X^T X = \begin{bmatrix} n & & & \\ & X_1^T X_1 & & \\ & & X_2^T X_2 & \\ & & & \ddots \\ & & & & X_k^T X_k \end{bmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$



$$H_0^1: \text{span}\{\mathbf{1}\} \quad H_A^1: \text{span}\{\mathbf{1}, W_1\}$$

$$H_0^2: \text{span}\{\mathbf{1}, W_2\} \quad H_A^2: \text{span}\{\mathbf{1}, W_1, W_2\}$$

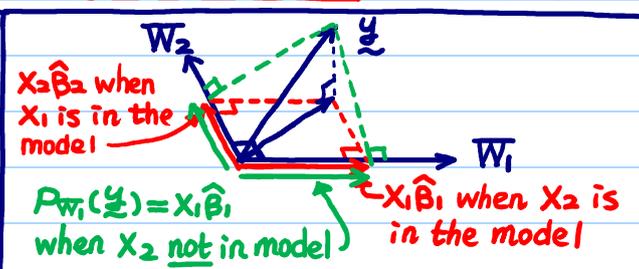
$H_0: y = \beta_0 \mathbf{1} + X_2 \beta_2 + X_5 \beta_5 + \dots + \epsilon$ (irrelevant)

$H_A: y = \beta_0 \mathbf{1} + X_2 \beta_2 + X_5 \beta_5 + \dots + X_1 \beta_1 + \epsilon$

$$\|y - \hat{y}\|^2 = \|P_{\Omega}(y)\|^2 + \|P_{\Omega^\perp}(y)\|^2 \quad (LN p. 25) \quad RSS_{H_0} - RSS_{H_A} = \|P_{W_1}(y)\|^2$$

$$= \|P_{W_1}(y)\|^2 + \|P_{W_2}(y)\|^2 + \dots + \|P_{W_k}(y)\|^2 + \|P_{\Omega^\perp}(y)\|^2$$

not true if not orthogonal



❖ Reading: Textbook, 1.4~1.6, 1.8

Some Properties of (Multivariate) Normal Distribution

(N1) Linear transformation of normal is still normal

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

$$E^*[(AZ+c)(AZ+c)^T] = E^*[(AZ)(AZ)^T] = E^*[AZZ^T A^T]$$

$$COV(Z) = \begin{bmatrix} COV(Z_1) & COV(Z_1, Z_2) \\ \tau \leftarrow & COV(Z_2) \end{bmatrix}$$

(N2) When 1st and 2nd moments are given, the normal distribution is specified. i.e., mean vector & variance-covariance matrix.

(N3) $Z = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} \sim$ normal, and Z_1, Z_2 uncorrelated (i.e., $COV(Z_1, Z_2) = 0$)

By (N3) $\Rightarrow Z_1, Z_2$ independent

$\underline{W} = \begin{bmatrix} W_1 \\ W_2 \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} Z$, $COV(W_1, W_2) = E^*(W_1 W_2^T) = E^*(A_1 Z Z^T A_2^T) = A_1 E^*(Z Z^T) A_2^T = A_1 \Sigma A_2^T$

(N4) $Z \sim N(\mu, \Sigma)$, $W_1 = A_1 Z$, $W_2 = A_2 Z$. Z can be generalized to k & $A_i \Sigma A_j^T = 0, i \neq j$

By (N3) (N4) $\Rightarrow W_1, W_2$ are independent iff $A_1 \Sigma A_2^T = 0$. If $\Sigma = \sigma^2 I$, then $A_1 \Sigma A_2^T = 0 \Leftrightarrow A_1 A_2^T = 0$

(N5) $Z \sim N(\mu, \Sigma)$, $W_1 = A_1 Z$, $W_2 = A_2 Z$, ..., $W_k = A_k Z$, and $COV(W_i, W_j) = 0$

length² of W_i for $1 \leq i < j \leq k$, $\Rightarrow W_1^T W_1, W_2^T W_2, \dots, W_k^T W_k$ are mutually independent.

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

$\Sigma^{-1} = (\Sigma^{1/2})^T (\Sigma^{1/2})$

- if Σ is non-singular, $(Z - \mu)^T \Sigma^{-1} (Z - \mu) \sim \chi_n^2$
- if Σ is singular and has rank $r (< n)$, \rightarrow The possible vectors of Z only occupy an r -dim subspace of \mathbb{R}^n

useful for the independence between sums of squares

standardization $\Sigma^{-1/2} (Z - \mu) \sim N(0, I)$

not unique $(Z - \mu)^T \Sigma^- (Z - \mu) \sim \chi_r^2$