



310.13
7874K
v.1

SAS의 통계처리방법(I)

1984. 11

경제기획원조사통계국

(527:0)



일 러 두 기

이 책은 SAS BASIC 에 관한 내용을 직원 교육용으로 활용하기 위하여 간략하게 간추린 것입니다. 통계국에서 1983년부터 설치운영하고 있는 SAS는 기본적인 통계분석 뿐만아니라 ETS (Econometion and Time Series), IMS/DL-1 (Information Management System/Data Language 1.) 과도 연결이 가능하며 Graphic OR (Operations Research) 등 다양한 분석기법을 취급할 수 있는 program package 인바 SAS를 이해하고자 하는 사람에게 도움이 되 고져 SAS 에 관한 최신정보를 입수하고 실무중심의 예제를 첨부 하여 실용적인 참고서를 만들려고 노력했습니다.

그러나 SAS가 가히 Software 의 혁명이라고 할만큼 기능이 다양하고 end user 에게 쉽게 사용할 수 있도록 만들어졌는데 반하여 이 책이 과연 입문서로서의 역할을 다할 수 있을런지는 자신 할 수 없으나 부족한 점은 계속하여 보완하고 개정하여 나가도록 노력할 계획이오니 양지하시기 바랍니다.

끝으로 이 책을 편찬하는데 심혈을 기울인 이경의처리관의 노고를 치하하는 바입니다.

자료처리과장

목 차

통계 측정 방법	9
(1) 명목 측정 (Nominal Measurement)	11
(2) 서열 측정 (Ordinal Measurement)	11
(3) 등간 측정 (Interval Measurement)	11
(4) 비율 측정 (Ratio Measurement)	12
1) 집중 경향의 지수	12
2) 분산도의 지수	13
3) SKEWNESS, KURTOSIS	16
CHAPTER 1. INTRODUCTION TO SAS	17
1. SAS PRODUCTS	19
(1) Statistical Analysis System	19
(2) Additional SAS Products: SAS/GRAPH	20
(3) Additional SAS Products: SAS/ETS	21
(4) Additional SAS Products: SAS/OR	22
(5) Biorhythm Chart	23
2. THE DATA STEP	24
(1) SAS Processing	24
(2) SAS Jobs	25

(3)	SAS Jobs	26
(4)	Documenting SAS Data Sets	27
(5)	The SAS Data Set	28
(6)	SAS Statements are Free-Format	29
(7)	A SAS Job with a DATA Step and Two PROC Steps	30
(8)	Overview of the DATA Step	31
(9)	Overview of the DATA Step	32
(10)	Overview of the DATA Step	33
(11)	SAS Job	34
3.	SAS PROCEDURE LIBRARY JCL	35
4.	SAS JCL	36
(1)	SAS View	37
(2)	A Simple SAS Job	38
(3)	A Listing of the Raw Data	39
5.	TEMPORARY DATA SETS	40
6.	A SIMPLE SAS JOB: OS BATCH	41
(1)	PROC PRINT Output: OS Batch	42
(2)	PROC MEANS Output: OS Batch	42
7.	A SIMPLE INTERACTIVE SAS JOB: TSO	43

8. SAS SYNTAX AND SAS DATA SETS	44
(1) Rules for Writing SAS Statements	44
(2) Structure of SAS Data Sets	45
(3) Naming SAS Data Sets and SAS Variables	46
(4) Procedure and Data Step	47
(5) Inputting Raw Data	48
(6) Formatted Input	49
(7) Selected Informats:	50
 CHAPTER 2. SAS DATA FILE	 51
1. CREATING VARIABLES AND EDITING VALUES	53
(1) Assignment Statements	53
(2) Example	54
(3) Types of Expressions	55
(4) SAS Functions	56
(5) Do and END Statements	57
(6) Sample	58
(7) Comments on Missing Values	58
(8) Selecting Observations	60
(9) The OUTPUT Statement	61
(10) Executing the DATA Step	62
(11) Sample	63
(12) Print the Data Set	63

(13)	Reading Selected Variables	64
(14)	The PUT Function in PROG Step	65
2.	CREATING FORMATS	67
(1)	The FORMAT Procedure	67
(2)	PROC FORMAT Options;	67
(3)	The VALUE Statement	68
(4)	Date, Time, and Datetime Informats and Formats	69
(5)	Date, Time, and Datetime Values	70
(6)	Concatenating SAS Data Sets	71
(7)	Interleaving SAS Data Sets	72
(8)	Match Merging	73
(9)	Match Merging	74
(10)	Match Merging	75
(11)	UPDATE Applying	76
(12)	UPDATE Applying	77
(13)	SAS Array	78
CHAPTER 3. 기초통계 처리		79
1.	CHART Procedure	81
(1)	기본형	81
(2)	일단형	82
(3)	Options	82

2. PLOT PROCEDURE	83
(1) 기본형	83
(2) 일반형	85
3. CORR Procedure	85
(1) Product-moment Couelation (Pearson)	85
(2) 기본형	86
(3) 일반형	87
(4) Options	87
4. FREQ Procedure (discrete Variable에 사용)	88
(1) 기본형	88
(2) 일반형	89
(3) TABLES Options	89
5. MEANS Procedure	92
(1) 기본형	92
(2) 일반형	92
(3) PROC MEANS Options;	93
(4) STATISTICS	94
6. UNIVARIATE Procedure	95
(1) 기본형	95
(2) STATISTICS	97

(3) 일 반 형	98
(4) PROC UNIVARIATE Options;	98
CHAPTER 4. 회귀분석과 분산분석에 관련된 기법	101
A. 회귀 분석 (Regression Analysis)	101
1. 회귀 분석	101
2. REG Procedure	101
3. STEPWISE Procedure	113
(1) PROC STEPWISE Option;	113
(2) PROC STEPWISE Option	113
(3) Sample	115
B. 분산 분석 (Analysis of Variance)	120
1. 분산 분석	120
2. ANOVA Procedure	121
3. Randomized Block Design (RBD)	124
C. t-검정 (t-Test)	127
1. Two Related Samples	127
2. Two Independent Samples	129
CHAPTER 5. 다변량 분석	131
A. DISCRIMINANT Analysis	133
1. Introduction	133
2. Related Procedures	133

3. Alternative Procedure	134
4. Background of DISCRIM Procedure	134
5. Outline of Use	136
 B. FACTOR Analysis	 136
1. Introduction	136
2. Related Procedures (in the Sense of Structural Analysis)	137
3. Background of FACTOR Procedure	138
4. Factor Rotation	139
5. Factor Extraction and Goodness of Fit Test on Factor Model	140
6. Outline of Use	140
 C. CLUSTER Analysis	 141
1. Introduction	141
2. Related Procedures	142
3. CLUSTER Procedures	143

통 계 측 정 방 법

통계 측정 방법

(1) 명목측정 (Nominal Measurement)

명목측정이라 함은 수로서의 특성을 전혀 갖지 않고 관찰대상의 고유한 속성을 분류하기 위해 숫자 등의 기호를 부여한 것이다. 예를 들면, 지역별 (1. 서울, 2. 부산, 3. 대구, 4. 인천), 성별 (1. 남, 2. 녀), 종교 (1. 기독교, 2. 불교, 3. 유교, 4. 이슬람교) 등을 측정할 경우다. 즉, 명목측정은 같다, 다르다. 이외에는 상대적 크기를 알 수 없음을 물론 덧셈, 곱셈과 같은 수학적 계산을 할 수 없는 것으로서 단순히 동일값을 표시하기 위한 수단의 역할만 하는 것이다. 따라서 숫자간의 크기나, 서열을 가정하는 통계적 분석 기법에서는 명목측정의 성격을 가진 자료는 이용될 수 없다.

(2) 서열측정 (Ordinal Measurement)

서열측정은 명목측정과 는 달리 같다, 다르다, 차이 뿐만 아니라 크기 서열까지를 제공하는 것으로서 예를들면, 직업분류에 있어서 1. 노무직, 2. 감독직, 3. 관리직을 각각 상, 중, 하로 분류하는 경우이다. 그러나 이 측정은 크다, 높다라는 사실 이외에는 어느정도의 차이가 나는가에 대한 정보는 제공하지 못한다. 따라서 서열측정 역시 수의 완전한 특성을 보유하지 못함으로 통계적 기법의 적용에는 한계가 있다고 하였다.

(3) 등간측정 (Interval Measurement)

등간측정이란 크다, 작다하는 서열적 성격 뿐만 아니라 얼마만큼 큰

가하는 차이 (distance)에 관한 정보를 갖는 측정이다. 예를 들면 20 세, 40 세, 60 세라는 세사람이 존재한다면 20 세와 40 세의 차이는 40 세와 60 세의 차이와 같다고 말할 수 있다. 그러나 60 세는 20 세보다 세배 늙었다고 말할 수 없다. 그 이유는 등간측정은 절대영점 (Absolute Zero Point)이 아닌 임의영점 (Arbitrary Zero Point)만 갖기 때문이다. 즉, 등간측정은 두 대상의 차이는 나타낼 수 있지만 비율에 대해서는 설명하지 못한다.

(4) 비율측정 (Ratio Measurement)

비율측정은 서열, 등간이외에 한 측정치는 다른 측정치의 두배 또는 세배 등의 비율에 관한 정보까지 제공한다. 즉, 비율측정은 절대영점에 정해져 있는 측정이다. 예를 들면, 무게 (길이)등에서 20 kg은 10 kg의 두배, 10 cm는 5 cm의 두배라 말할 수 있다.

1) 집중 경향의 지수

최빈도 (Mode), 중앙치 (Median), 평균 (Mean)

최빈도 (Mode)

최빈도란 한분포에서 빈도가 가장 높은 수치를 말한다. 이때 모든 변수가 똑같은 빈도를 갖는 분포가 있을 수 있다. 이러한 경우에 최빈도는 없다고 본다. 또한 가장 높은 빈도를 2개이상 갖는 분포도 있다. 이러한 분포를 상봉분포 (bimodal distribution) 삼봉분포 (trimodal distribution), 다봉분포 (multimodal distribution)을 이룬다.

중앙치 (Median)

중앙치는 역시 집중경향의 지수의 하나로서 한 분포안에 포함된 전체사례 (N)을 양등분하는 점에 해당하는 수치를 말한다. 즉, 점수들이 순서대로 나열되었을 때 제일 가운데 있는 점수가 중앙치가 된다.

평균 (Mean)

평균은 모든 점수를 다합한 값을 전체 사례수로 나눈 값이다. 평균치는 가장 민감한 집중경향의 지수로서 대부분의 통계절차에 밀접한 관계가 있다. 여기서 평균이 가지고 있는 몇가지 속성은 다음과 같다.

- ① 각 사례의 점수로부터 평균값을 뺀 즉 편차의 합은 0이다.
즉 $\sum (X - M) = 0$
- ② 평균은 특히 분포가 Skewness 되어 있을때 극단수치 (extremes core)에 매우 민감하다.
즉 평균은 각 수치들의 무게의 중심으로서 지렛대의 균형점 구실을 하는데 한두개의 수치라도 평균으로부터 상대적으로 먼지점에 걸려 있을때 평균은 그 분포의 대표적 수치로서 기능을 다할 수 없는 것이다.
- ③ 평균치로부터의 편차의 자승의 합 즉 $\sum (X - M)^2$ 은 다른 어떤 점수들을 기준으로 해서 얻어진 편차의 제곱들의 총합보다 항상 적다. 평균치는 편차의 자승의 값이 최소가 되는 집중 경향치로서 최소자승의 합을 구함으로써 평균치를 알아낼 수 있다. 이러한 방법이 최소자승법 (Least Squares Method)이다.

2) 분산도의 지수

범위 (Range), 분산 (Variance), 표준편차 (Standard Deviation)

표준점수 (Standard Score)

분산도 (dispersion)란 한 분포안의 사례들이 집중경향치를 중심으로 얼마나 밀집, 혹은 분산되어 있는가 하는 정도를 말한다. 분산도는 집중경향치와 함께 분포를 해석 기술하는데 중요한 역할을 하며 지수를 통해 측정되는데 가장 빈번히 사용되는 지수들은 범위, 분산, 표준편차 등이다.

범위 (Range)

범위는 한 분포내에서 최고치와 최저치 사이의 간격으로 최고치에서 최저치를 뺀 값을 말하며 보통 R로 표시한다. 즉 $R = \text{Max} - \text{Min}$ 이다.

분산 (Variance)

분산 또는 변량이란 편차를 자승하여 모두 합한 후에 이 값을 (진제사예수 - 1)로 나눈 값이다. 한 분포안의 값들이 평균으로부터 떨어져 있는 방향이 편차 (variance)이다. 이들 편차를 모두 합하면 $\sum d = 0$ 이 된다. 이들 편차의 총합이 0이 되지 않도록 하기 위하여 편차를 자승한 값이 바로 분산이다.

$$S^2 = \frac{\sum (X - M)^2}{N - 1}$$

표본의 변량 또는 표준편차를 계산할 때 분모는 보통 $N - 1$ 이 되는데 이것은 모집단의 변량과 표준편차에 대한 불편 추정치 (unbiased estimated)가 되기 때문이다. 사례수 (N)이 작을 때에는

불편수정치 $N-1$ 을 사용하지만 N 이 클 때에는 $N-1$ 대신 N 을 사용해도 별차이는 없다. 변량을 표시하는 방법으로 S^2 과 σ^2 이 사용되는데 통상 S^2 은 표본의 변량을 말하고 σ^2 은 모집단의 변량을 말한다.

표준편차 (Standard Deviation)

분산도의 지수로서 변량보다 더 일반적으로 사용되는 것이 표준편차이다. 이는 변량에 제곱근을 씌운 이유는 변량을 구하기 위하여 편차를 자승하여 합쳤으므로 그 반대로 다시 제곱근을 씌운 것이다. 표준편차를 나타내는 기호로 모집단의 경우에 σ , 표본의 경우 S 를 사용하는 것이 보통이다.

표준점수 (Standard Score)

편차, 분산, 표준편차들은 어떤 점수들이 한 분포내에서 차지하는 상대적 위치를 파악하고자 하는 경우였으나 분포가 다른 두 점수들은 비교해 보고자하는 경우가 있다. 코끼리무게와 사람의 몸무게와 같이 평균치와 표준편차가 다른 두 분포내의 어떤 두 점수들을 비교하고자 하는 경우에 있다. 이와 같이 한 분포내에서 차지하는 상대적 위치나 분포가 다른 두 점수들을 직접 비교하는 경우에 사용되는데 이를 표준점수 또는 Z score라고도 한다.

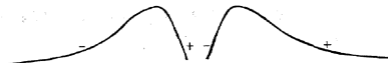
$$\text{일반공식은 } Z = \frac{X - M}{S}$$

X = 원점수, M = 평균, S = 표준편차

3) SKEWNESS, KURTOSIS

負의偏布(negative skew)

正의偏布(positive skew)



Skewness

어떤 분포가 정상분포에서의 대칭성 (symmetry)을 어느 정도 만족시키는가를 나타내는데 이 역시 정상분포일때 0의 값을 갖는다. skewness가 0보다 크다 ($S > 0$)는 것은 다수의 값이 평균보다 낮은곳에 집중하고 극소수의 값들이 평균보다 높은곳에 위치하고 있어 positive skewness를 이룬다. skewness가 0보다 작을 때 ($S < 0$)에는 이와 정반대로 해석하면 된다.

正常分布(bell-shaped)

뾰족한 分布(peaked)

납작한 分布(flat)



Kurtosis

분포가 어느 정도 뾰족하게 솟아 있는지를 표시하는 것으로서 정상 분포일 때 0의 값을 갖는다. kurtosis가 0보다 클때($K > 0$) 분포는 정상 분포보다 뾰족한 (peaked) 모양을 나타내고 0보다 작을때 ($K < 0$)는 평평한 (flat) 모양을 나타내게 된다.

CHAPTER 1. INTRODUCTION TO SAS

CHAPTER 1. INTRODUCTION TO SAS

1. SAS PRODUCTS

(1) Statistical Analysis System

What is SAS?

SAS is a computer software system that consists of several products that provide tools for data entry, data management, and data analysis:

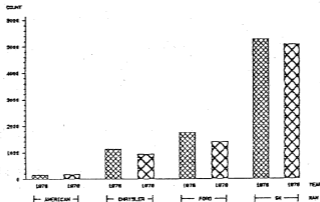
RETRIEVAL	flexible input techniques.
TRANSFORMATIONS	programming language with statistical and mathematical functions.
MAINTENANCE	storing, documenting, updating, and editing.
MANIPULATION	sorting, subsetting, concatenation, and merging.
REPORT WRITING	printing information using program statements.
PRINTER GRAPHICS	charts and two-dimensional plots.
DATA REDUCTION AND SUMMARIZATION	descriptive statistics.
STATISTICAL ANALYSIS	from simple crosstabulations to complex multivariate techniques.

(2) Additional SAS Products: SAS/GRAPH

A device-intelligent color graphics system that features:

- o charts
- o plots
- o maps
- o slides
- o reports

U.S. Passenger Car Production

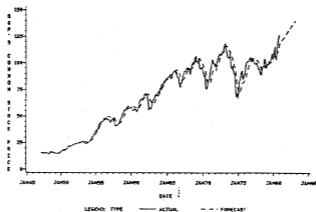


(3) Additional SAS Products: SAS/ETS

An econometric and time series analysis system that features:

- o time series forecasting
- o time series regression techniques
- o simultaneous equation modeling techniques
- o financial reports

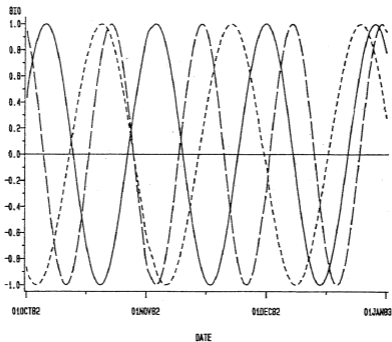
PROC FORECAST Output
Displayed with PROC Gplot



(5) Biorhythm Chart

Birthday Sept. 10, 1952.

BIORHYTHM CHART

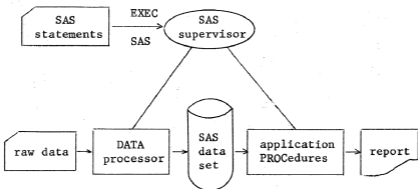


LEGEND: PHYTHM — EMOTION INTELLECT ——— PHYSICAL

2. THE DATA STEP

(1) SAS Processing

SAS consists of a data-handling language and a library of procedures that work together as a system.



- o A supervisor program, written in re-entrant Assembler, directs the execution.
- o All storage is allocated as needed.
- o The supervisor links to the programs on the library.
- o Each procedure is one or two separate load modules on the library.

(2) SAS Jobs

All SAS jobs are a sequence of SAS steps.

There are only two kinds of SAS steps:

- o DATA steps prepare SAS data sets
- o PROCedure steps analyze or process SAS data sets.

Example:

Given : A company has recored revenue and expense data yearly from 1978 through 1981.

Objective : Compute the difference between expenses and revenue (income) for each year and compute the average income across all the years.

data fields →

year	revenue	expenses
78	45000	32000
79	53000	36000
80	54000	38000
81	65000	43000

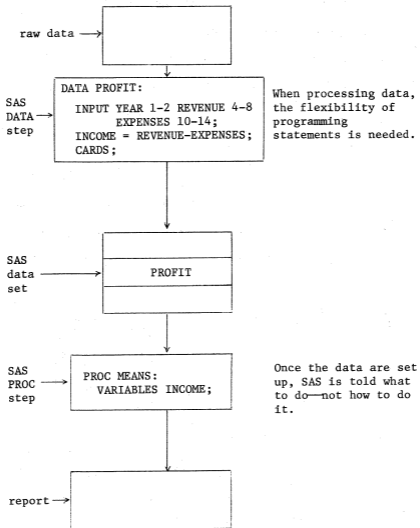
raw data →

variables →

YEAR	REVENUE	EXPENSES	INCOME= REVENUE-EXPENSES
78	45000	32000	13000
79	53000	36000	17000
80	54000	38000	16000
81	65000	43000	22000

SAS data set →

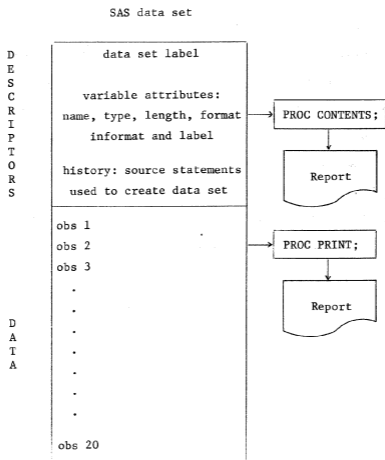
(3) SAS Jobs



(4) Documenting SAS Data Sets

A SAS data set contains:

- o descriptor records
- o data records.



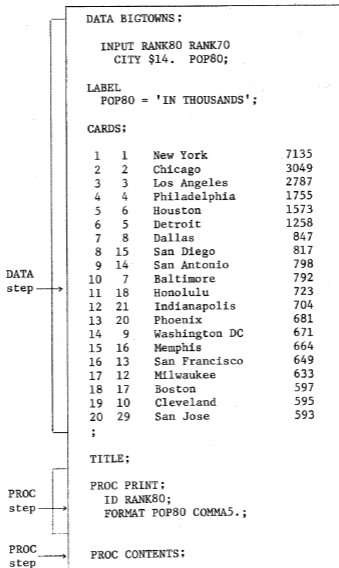
(5) THE SAS DATA SET

- o DATA steps create SAS data sets.
- o SAS data sets are rectangular or flat.
- o You refer to a SAS data set by its name. You refer to SAS variables by their names.
- o Every SAS data set has a name and is physically stored on disk or tape. In simple jobs, the SAS data sets are stored on temporary space, but they can be stored permanently.
- o Data must be in the form of a SAS data set before they can be analyzed by SAS procedures.
- o There is no limit to the number of DATA steps that can be included in a single job.
- o If a raw file is to be produced, the DATA step can create a null SAS data set, saving time and space.

(6) SAS Statements are Free-Format

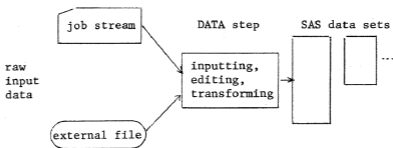
- o All SAS statements end with a semicolon (;).
- o Statements begin and end anywhere.
- o Several statements can be on the same line.
- o One statement can continue over several lines.
- o SAS words are separated by one or more blanks.

(7) A SAS Job with a DATA Step and Two PROC Steps.



(8) Overview of the DATA Step

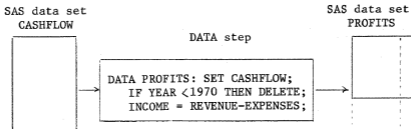
Retrieval



Transferring with reshaping

Given : SAS data set CASHFLOW contains the variables
YEAR, REVENUE, and EXPENSES.

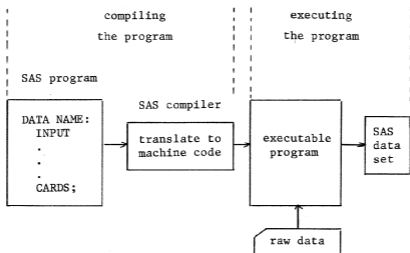
Objective: Delete observations for years before 1970
and compute the income for each remaining
year.



(10) Overview of the DATA Step

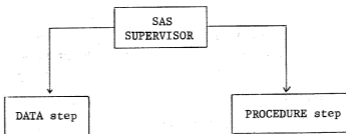
The DATA step involves writing a compact SAS program to process data. SAS processes the program in two steps.

- o The program is compiled.
- o The program is executed.



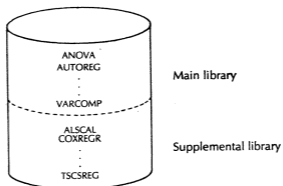
(11) SAS Job

- o DATA steps prepare SAS data sets.
- o PROCEDURE (or PROC) steps process SAS data sets.

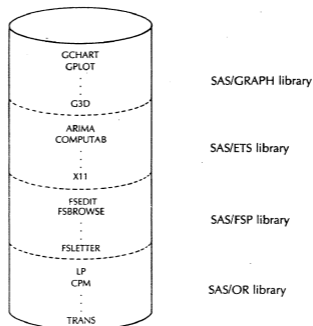


3. SAS PROCEDURE LIBRARY JCL

- Procedures are programs designed to process SAS data sets.
- A procedure is called by name from the SAS procedure library.



SAS offers optional 'add-on' procedure libraries



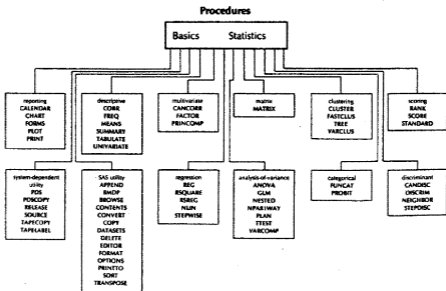
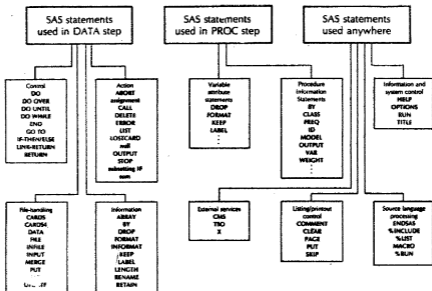
4. SAS JCL

```
//SASSAS30 JOB  
//SAS EXEC SAS  
//SYSIN DD *
```

SAS PROGRAM

```
/*  
//
```

(1) SAS VIEW



(2) A simple SAS job

Example :

Given : The name, sex, age, height, and weight of each student in a summer camp were recorded and stored with the following record format.

DATA FIELD	FIELD
DESCRIPTION	POSITION
NAME	1-10
SEX	11
AGE	13-14
HEIGHT	16-19
WEIGHT	21-25

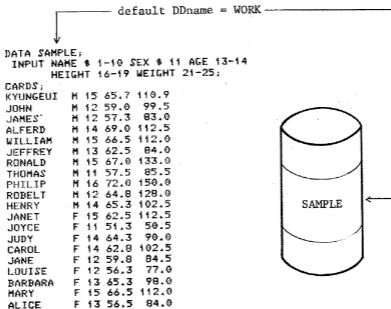
- Objectives :
1. Obtain a listing of the data.
 2. Compute some descriptive statistics for the heights and weights such as the mean, minimum, and maximum.
 3. Produce a plot of height versus weight.

(3) A Listing of the Raw Data

NAME							SEX	AGE	HEIGHT	WEIGHT																			
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
K	Y	U	N	G	E	K	1	M	15	63	7	110	9																
J	O	H	N					M	12	59	0	99	5																
J	A	M	E	S				M	12	57	3	82	0																
A	L	F	R	E	D			M	14	69	0	112	5																
W	I	L	L	I	A	M		M	15	66	5	112	0																
J	E	F	F	R	E	V		M	13	62	5	84	0																
R	O	M	A	L	D			M	15	67	0	133	0																
T	H	O	M	A	S			M	11	57	5	85	0																
P	H	I	L	I	P			M	16	72	0	150	0																
R	O	B	E	R	T			M	12	64	8	128	0																
H	E	N	R	I				M	14	63	5	102	5																
J	A	M	E	T				F	15	62	5	112	5																
J	O	I	C	E				F	11	51	3	50	5																
J	U	D	I					F	14	64	3	90	0																
C	A	R	O	L				F	14	62	8	102	5																
J	A	M	E					F	12	59	8	84	5																
L	O	U	I	S	E			F	12	56	3	77	0																
B	A	R	B	A	R	A		F	13	65	3	98	0																
M	A	R	K					F	15	66	5	112	0																
A	L	I	C	E				F	13	56	5	84	0																

Note: These data can be on cards, disk, tape, or lines entered on a terminal.

5 Temporary Data Sets



The note on the SAS log for this DATA step would read:

NOTE: DATA SET WORK. SAMPLE HAS 10 OBSERVATIONS AND 5
VARIABLES. 504 OBS/TRK.

The data set WORK. SAMPLE has been stored in work space
and will be deleted at the end of the job or session.

The data set may be referred to as SAMPLE or WORK. SAMPLE

```
PROC PRINT DATA=WORK. SAMPLE;
```

is equivalent to

```
PROC PRINT DATA=SAMPLE;
```

6. A Simple SAS Job: OS Batch

```
//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
DATA SAMPLE;
  INPUT NAME $ 1-10 SEX $ 11 AGE 13-14
        HEIGHT 16-19 WEIGHT 21-25;
CARDS;
KYUNGEUI M 15 65.7 110.9
JOHN M 12 59.0 99.5
JAMES M 12 57.3 83.0
ALFERD M 14 69.0 112.5
WILLIAM M 15 66.5 112.0
JEFFREY M 13 62.5 84.0
RONALD M 15 67.0 133.0
THOMAS M 11 57.5 85.5
PHILIP M 16 72.0 150.0
ROBELT M 12 64.8 128.0
HENRY M 14 65.3 102.5
JANET F 15 62.5 112.5
JOYCE F 11 51.3 50.5
JUDY F 14 64.3 90.0
CAROL F 14 62.8 102.5
JANE F 12 59.8 84.5
LOUISE F 12 56.3 77.0
BARBARA F 13 65.3 98.0
MARY F 15 66.5 112.0
ALICE F 13 56.5 84.0
;
PROC PRINT DATA=SAMPLE;
PROC MEANS DATA=SAMPLE;
  VAR WEIGHT HEIGHT;
PROC PLOT DATA=SAMPLE;
  PLOT WEIGHT*HEIGHT;
/*
//
```


(1) PROC PRINT output: OS batch

14:47 WEDNESDAY, OCTOBER 31, 1984 1

SAS

OBS	NAME	SEX	AGE	HEIGHT	WEIGHT
1	KYUNGEUJI	M	15	65.7	110.9
2	JOHN	M	12	59.0	99.5
3	JAMES	M	12	57.3	83.0
4	ALFERD	M	14	69.0	112.5
5	WILLIAM	M	15	66.5	112.0
6	JEFFREY	M	13	62.5	84.0
7	RONALD	M	15	67.0	133.0
8	THOMAS	M	11	57.5	85.5
9	PHILIP	M	16	72.0	150.0
10	ROBELT	M	12	64.3	123.0
11	HENRY	M	14	65.3	102.5
12	JANET	F	15	62.5	112.5
13	JOYCE	F	11	51.3	50.5
14	JUDY	F	14	64.3	90.0
15	CAROL	F	14	62.8	102.5
16	JANE	F	12	59.3	84.5
17	LOUISE	F	12	56.3	77.0
18	BARBARA	F	13	65.3	98.0
19	MARY	F	15	66.5	112.0
20	ALICE	F	13	56.5	84.0

(2) PROC MEANS output: OS batch

VARIABLE	N	MEAN	STANDARD DEVIATION	MINIMUM VALUE	MAXIMUM VALUE	STD DEVIATION OF MEAN	SUM	VARIANCE	L.F.
WEIGHT	20	106.5450000	27.78137704	56.50000000	150.00000000	4.49326478	2031.90000000	496.43964737	12.150
HEIGHT	20	62.54500000	5.07174601	51.30000000	72.00000000	1.17603760	1251.90000000	29.33966312	8.110

7. A Simple Interactive SAS Job: TSO

```

READY
SAS
NOTE: SAS RELEASE 82.2
      AT SAS INSTITUTE INC. (0) (00000).
1?
DATA CLASS;
2?
  INPUT NAME $ 1-8 SEX $ 11 AGE 13-14
3?      HEIGHT 16-19 WEIGHT 21-25;
4?
CARDS;
5>
KYUNGEUI M 15 65.7 110.9
6>
JONH     M 12 59.0 99.5
7>
JAMES   M 12 57.3 83.0
8>
ALERED  M 14 69.0 112.5
9>
WILLEAM M 15 66.5 112.0
10>
JEEEREY M 13 62.5 84.0
11>
RONALD  M 15 67.0 133.0
12>
THOMAS  M 11 57.5 85.0
13>
PHILIP  M 16 72.0 150.0
14>
ROBERT  M 12 64.8 128.0
15>
HENRY   M 14 63.5 128.0
16>
JAMET   E 15 62.5 112.5
17>
JOYCE   E 11 51.3 50.5
18>
JUDY    E 14 64.3 90.0
19>
CAROL   E 14 62.8 102.5
20>
JANE    E 12 59.8 84.5
21>
LOUISE  E 12 56.3 77.0
22>
BARBARA E 13 65.3 98.0
23>
MARY    E 15 66.5 112.0
24>
ALICE   E 13 56.5 84.0
25>
RUN;
NOTE: DATA SET WORK CLASS HAS 19 OBSERVATIONS

```

8. SAS Syntax and SAS Data Sets

(1) Rules for writing SAS statements

SAS statements begin with an identifying keyword and end with a semicolon.

```
DATA CLASS:
```

```
    INPUT NAME $1-10 SEX $11 AGE 13-14
```

```
        HEIGHT 16-19 WEIGHT 21-25;
```

```
    CARDS;
```

data lines

```
PROC PRINT DATA=CLASS;
```

```
PROC MEANS DATA=CLASS;
```

```
    VARIABLES HEIGHT WEIGHT;
```

```
PROC PLOT DATA=CLASS;
```

```
    PLOT WEIGHT*HEIGHT;
```

SAS statements are free-format.

- o They can begin anywhere, end anywhere
- o One statement can continue over several lines.
- o Several statements can be on a line.
- o Blanks-as many as you want-separate fields. Special characters also separate fields.

Note: We recommend that DATA and PROC statements start in column 1 and that other statements be indented.

(2) Structure of SAS Data Sets

A SAS data set is a collection of data values arranged in a rectangular table.

VARIABLES

	NAME	SEX	AGE	HEIGHT	WEIGHT
observation 1	KYUNGEUI	M	15	65.7	110.9
observation 2	JOHN	M	12	59.0	99.5
observation 3	JAMES	M	12	57.3	83.0
.
observation 20	ALICE	F	13	56.5	84.0

The columns in the table are called variables.

- o Variables correspond to fields of data.
- o Each variable is given a name.
- o There are two kinds of variables.

character variables: NAME, SEX (1-200 long)

numeric variables: AGE, HEIGHT, WEIGHT (floating)

The rows in the table are called observations (or records).

There is no limit on the number of observations.

Note: The rectangular structure of a SAS data set implies that every variable must exist for each observation.

(3) Naming SAS Data Sets and SAS Variables

Rules for names:

- o 1-8 characters
- o start with A-Z or _
- o continue with numbers, letters, or underscores.

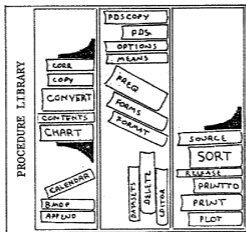
Suggestion: Choose meaningful data set and variable names.

Note: Under CMS SAS, the underscore, '_', is not allowed in SAS data set names.

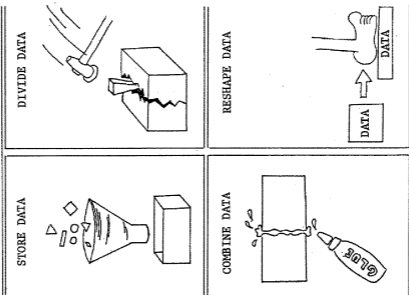
Example:

```
DATA CLASS;  
  INPUR NAME $ 1-10 SEX $11 AGE 13-14  
  HEIGHT 16-19 WEIGHT 21-25;  
CARDS;
```

(4) Procedure and Data Step.



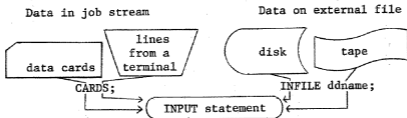
DATA STEP: DATA HANDLING



- The SAS Procedure Library contains PROCs or procedures to process SAS data sets.
- Each is accessed by name from the SAS procedure library.

(5) Inputting Raw Data

SAS can handle data in virtually any form, from almost any input file.



Note: Regardless of where the data are stored, the same INPUT statement is used.

Functions of the CARDS, INFILE and INPUT statements

- o Reading the data.

INPUT statement reads raw data lines

- assigns names to the SAS variables that correspond to the fields
- has three modes: COLUMN, LIST, and FORMATTED.

- o Pointing to the data file.

CARDS statement indicates to SAS that data records follow immediately

INFILE statement points SAS to an external file where the raw data are stored.

(6) Formatted Input

Specify the starting location and field widths (similar to FORTRAN and PL/I). Move an input "pointer" to the starting position of the field, then specify the variable and an informat.

INPUT pointer control variable informat;

informats: W. numeric width
 W.d numeric with decimal
 \$W. character

pointer controls: @n go to column n
 +n move the pointer n positions
 W. informats advance the pointer

INPUT NAME \$8. @11 SEX \$1. +1 AGE 2.

+1 HEIGHT 4. +1 WEIGHT 5.;

NAME										SEX	AGE	HEIGHT	WEIGHT																	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	
K	Y	U	N	S	E	K	E			M	1	5	6	5	.	7				1	1	0	.	9						
J	O	H	N							M	1	2	5	9	.	0				9	9	.	5							
J	A	M	E	S						M	1	2	5	7	.	3				8	3	.	0							
A	L	F	R	E	D					M	1	4	6	9	.	0				1	1	2	.	5						

With formatted input you can read data in nonstandard numeric or character formats.

(7) Selected Informats:

w.	standard numeric
w.d	standard numeric with decimal
\$w.	standard character
\$CHARw.	characters with blanks
HEXw.	numeric hexadecimal
\$HEXw.	character hexadecimal
IBw.d	integer binary
PIBw.d	positive integer binary
PDw.d	packed decimal
PKw.	unsigned packed decimal
RBw.d	real binary (floating point)
ZDw.d	zoned decimal
ZDBw.d	zoned decimal with blanks
CBw.d	
\$CBw.	column binary
PUNCHd.	
ROWw.d	
COMMAw.d	commas in numbers
Ew.	scientific notation
EZw.d	blanks are zeros
\$VARYINGw.	varying-length character values

See pages 388-398 in the SAS User's Guide: Basics, 1982 Edition for additional information on SAS informats.

CHAPTER 2. SAS DATA FILE 에 관하여

CHAPTER 2. SAS DATA FILE에 관하여

1. CREATING VARIABLES AND EDITING VALUES

(1) Assignment Statements

Assignment statements are used to create new variables and to modify values of existing variables. SAS evaluates an expression and assigns the result to a variable.

```
variable=expression;
```

(2) Example

- o Read three variable (YEAR, REVENUE, and EXPENSES) into a SAS data set.
- o Add a variable named INCOME, which is the difference between REVENUE and EXPENSES.
- o Change the values of YEAR from 2 digits to 4 digits.

```
//SPSSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
DATA PROFITS;
  INPUT YEAR REVENUE EXPENSES;
  INCOME=REVENUE-EXPENSES;
  YEAR=YEAR+1900;
CARDS;
00 5650 1050
01 6280 1140
;
PROC PRINT;
/*
//
```

```
SAS
OBS YEAR REVENUE EXPENSES INCOME
1 1980 5650 1050 4600
2 1981 6280 1140 5140
```

Program data vector

YEAR	REVENUE	EXPENSES	INCOME

Note: Any variable defined by an assignment statement is included in the program data vector.

(3) Types of Expressions

- o simple arithmetic operation: + - * / **
 - X2=X; move the value
 - SUM=X+Y addition
 - DIF=X-Y subtraction
 - TWICE=X*2; multiplication
 - HALF=X/2; division
 - CUBIC=X**3; exponentiation
 - Y=-X; change the sign
- o constants
 - N=0; numeric constant
 - SEX='FEMALE'; character constant
- o complex expressions
 - priority of evaluation () ** * / + -
 - A=X+Y+Z; left to right
 - A=X+Y*Z; operator precedence
 - A=X/Y/Z; left to right
 - A=X/(Y/Z); parenthetical
- o functions
 - variable=FUNCTIONNAME (argument1, argument2,...);
 - S=SQRT(X);
 - A=ABS(X);
 - Z=ABS (SQRT(X)-2);

(4) SAS Functions

- o Selected functions that compute simple statistics.

SUM	sum
MEAN	arithmetic mean
VAR	variance
MIN	minimum value
MAX	maximum value
STD	standard deviation

Example:

Given : Temperature data at a specific location are recorded every hour on the hour for several days. Each record in a file represents one day and contains the date and the 24 recorded temperatures for that date.

Objective: Create a SAS data set that contains the date, the 24 hourly temperatures, the average temperature, the minimum temperature and the maximum temperature for each day.

```
DATA TEMP;
  INPUT DATE $ 1-7 @11 (T1-24) (2.);
  AVGTEMP=MEAN (OF T1-T24);
  MINTEMP=MIN (OF T1-T24);
  MAXTEMP=MAX (OF T1-T24);
  CARDS:
data lines
```

program data vector

DATE	T1	...	AVGTEMP	MINTEMP	MAXTEMP

(5) DO and END Statements

Example : Execute several statements when a condition
is met.

Given : Salary information for a group of employees.

Objective : If the department number is 201, then define
the department name to be SALES and set
gross pay equal to salary plus commission.
Otherwise, the department name is ADMIN and
the gross pay is simply equal to salary.

(6) Sample

```
//SPSSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
DATA EMPLOY;
  INPUT NAME# 1-8 DEPTNO 10-12
         COM 14-17 SALARY 19-23;
  IF DEPTNO=201 THEN
    DO;
      DEPT='SALES';
      GROSSPAY=COM+SALARY;
    END;
  ELSE
    DO;
      DEPT='ADMIN';
      GROSSPAY=SALARY;
    END;
  CARDS;
  JOHNSON 201 1500 18000
  MOSSER 101 21000
  LARKIN 101 24000
  GARRETT 201 4800 18000
  PROC PRINT;
```

```
              SAS
OBS      NAME      DEPTNO      COM      SALARY      DEPT      GROSSPAY
  1     JOHNSON      201        1500      18000      SALES      19500
  2     MOSSER       101          .      21000      ADMIN      21000
  3     LARKIN       101          .      24000      ADMIN      24000
  4     GARRETT      201        4800      18000      SALES      22800
```

(7) Comments on Missing Values

Missing values propagate through arithmetic expressions.

```
DATA;
  INPUT CHECKING SAVINGS;
  TOTAL1=CHECKING+SAVINGS;
  TOTAL2=SUM (CHECKING, SAVINGS);
  CARDS;
```


100 2000

300

CHECKING	SAVINGS	TOTAL1	TOTAL2
100	2000	2100	2100
300	.	.	300

Missing values compare as minus infinity.

DATA;

INPUT PAYMENT DUE;

IF PAYMENT < DUE THEN STATUS='PAST DUE';

ELSE STATUS='PAID';

CARDS;

10 20

25 25

10

PAYMENT	DUE	STATUS
10	20	PAST DUE
25	25	PAID
.	10	PAST DUE

(8) Selecting Observations

The subsetting IF statement

General form of the subsetting IF statement:

IF expression;

The subsetting IF statement is equivalent to:

IF \neg (expression) THEN DELETE;

The subsetting If statement tells SAS which observations to include in the output SAS data set. The statement works like a gate; it allows an observation to pass when the expression is true.

Example: Take a subset of the data.

```
DATA HISTORY;
```

```
INPUT YEAR REVENUE;
```

```
IF YEAR 1970 THEN DELETE;←
```

```
CARDS;
```

```
DATA HISTORY;
```

```
INPUT YEAR REVENUE;
```

```
IF YEAR =1970;←
```

```
CARDS;
```

(9) The OUTPUT Statement

```
DATA HISCHOOL COLLEGE;
    INPUT NAME $ 1-8 SEX $ 10 YRS-EDUC 12-13;
    IF YRS_EDUC <= 12 THEN OUTPUT HISCHOOL;
    IF YRS_EDUC > 12 THEN OUTPUT COLLEGE;
    CARDS;
KATHRYN  F 16
GEORGE   M 12
WILLIAM  M 18
JENNIFER F 12
CYNTHIA  F 16
```

input buffer

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	...	80

	NAME	SEX	YRS_EDUC
program data vector	<input type="text"/>	<input type="text"/>	<input type="text"/>

SAS data set
HISCHOOL

NAME	SEX	YRS_EDUC
GEORGE	M	12
JENNIFER	F	12

SAS data set
COLLEGE

NAME	SEX	YRS_EDUC
KATHRYN	F	16
WILLIAM	M	18
CYNTHIA	F	16

(10) Executing the DATA Step

```
DATA NEWSMP;  
  SET SAMPLE;  
  BIRTH-YR=1984-AGE;
```

SAS data set

SAMPLE

NAME	SEX	AGE	HEIGHT	WEIGHT
KYUNGEUI	M	15	65.7	110.9
JOHN	M	12	59.0	99.5
JAMES	M	12	57.3	83.0
ALFRED	M	14	69.0	112.5

program
data
vector

NAME	SEX	AGE	HEIGHT	WEIGHT

SAS data set

NEWSAMP

NAME	SEX	AGE	HEIGHT	WEIGHT	BIRTH_YR
KYUNGEUI	M	15	65.7	110.9	1969
JOHN	M	12	59.0	99.5	1972
JAMES	M	12	57.3	83.0	1972
ALFRED	M	14	69.0	112.5	1970

(11) Sample

```

//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SYSIN DD *
DATA SAMPLE;
INPUT NAME $ 1-10 SEX $ 11 AGE 13-14
HEIGHT 16-19 WEIGHT 21-25;
CARDS;
KYUNGEHI M 15 65.7 110.9
JOHN M 12 59.0 99.5
JAMES M 12 57.3 83.0
ALFRED M 14 69.0 112.5
WILLIAM M 15 66.5 112.0
JEFFREY M 13 62.5 84.0
RONALD M 15 67.0 133.0
THOMAS M 11 57.5 85.5
PHILIP M 16 72.0 150.0
ROBELT M 12 64.8 128.0
HENRY M 14 65.3 102.5
JANET F 15 62.5 112.5
JOYCE F 11 51.3 50.5
JUDY F 14 64.3 90.0
CAROL F 14 62.8 102.5
JANE F 12 59.8 84.5
LOUISE F 12 56.3 77.0
BARBARA F 13 65.3 98.0
MARY F 15 66.5 112.0
ALICE F 13 56.5 84.0
;
DATA NEWSAMP;
SET SAMPLE;
BIRTH_YR=1984-AGE;
PROC PRINT DATA=NEWSAMP;
/*
//

```

(12) Print the Data Set

```

SAS

```

PRS	NAME	SEX	AGE	HEIGHT	WEIGHT	BIRTH_YR
1	KYUNGEHI	M	15	65.7	110.9	1969
2	JOHN	M	12	59.0	99.5	1972
3	JAMES	M	12	57.3	83.0	1972
4	ALFRED	M	14	69.0	112.5	1970
5	WILLIAM	M	15	66.5	112.0	1969
6	JEFFREY	M	13	62.5	84.0	1971
7	RONALD	M	15	67.0	133.0	1967
8	THOMAS	M	11	57.5	85.5	1978
9	PHILIP	M	16	72.0	150.0	1968
10	ROBELT	M	12	64.8	128.0	1972
11	HENRY	M	14	65.3	102.5	1970
12	JANET	F	15	62.5	112.5	1969
13	JOYCE	F	11	51.3	50.5	1973
14	JUDY	F	14	64.3	90.0	1970
15	CAROL	F	14	62.8	102.5	1970
16	JANE	F	12	59.8	84.5	1972
17	LOUISE	F	12	56.3	77.0	1972
18	BARBARA	F	13	65.3	98.0	1971
19	MARY	F	15	66.5	112.0	1969
20	ALICE	F	13	56.5	84.0	1971

(13) Reading Selected Variables

Example: DROP and KEEP input data set options.

```
DATA SUBSET;
```

```
SET CLASS (DROP=HEIGHT WEIGHT);←
```

```
DATA SUBSET;
```

```
SET CLASS (KEEP=NAME SEX AGE);←
```

SAS data set

CLASS

NAME	SEX	AGE	HEIGHT	WEIGHT

NAME	SEX	AGE

program
data
vector

SAS data set
SUBSET

NAME	SEX	AGE

(14) The PUT Function in PROG Step

Example:

```
//SPCSAS30 JOB CLASS=V
//SAS      EXEC SAS
//SYSIN   DD *
DATA EXAMPLE;
  INPUT DEC @@;
  CNUM = PUT(DEC,4.);
  HEX  = PUT(DEC,HEX.);
  ROMAN= PUT(DEC,ROMAN.);
  WORDS= PUT(DEC,WORDS15.);
CARDS;
1 2 3 4 5 6 7 8 9 10
50 100 1000
;
PROC PRINT ;
  VAR CNUM HEX ROMAN WORDS DEC;
  ID DEC;
  TITLE PUT FUNCTION;
PROC CONTENTS;
RUN;
/*
//
```

PUT FUNCTION

DEC	CNUM	HEX	ROMAN	WORDS	DEC
1	1	00000001	I	ONE	1
2	2	00000002	II	TWO	2
3	3	00000003	III	THREE	3
4	4	00000004	IV	FOUR	4
5	5	00000005	V	FIVE	5
6	6	00000006	VI	SIX	6
7	7	00000007	VII	SEVEN	7
8	8	00000008	VIII	EIGHT	8
9	9	00000009	IX	NINE	9
10	10	0000000A	X	TEN	10
50	50	00000032	L	FIFTY	50
100	100	00000064	C	ONE HUNDRED	100
1000	1000	000003E8	M	ONE THOUSAND	1000

PUT FUNCTION

17:39 THURSDAY, NOVEMBER 11, 1964 2

CONTENTS OF SAS DATA SET WORK.EXAMPLE

TRACKS USED=1 SUBRENTS=2 OBSERVATIONS=13 CARRIED BY DS JOB SPC6AS10 ON CRUID 00-3021-000000

AT 17:39 THURSDAY, NOVEMBER 11, 1964 BY SAS RELEASE 67.3 USASAP=SYSA930,1170A02,AF103,SPC6AS10,00000001 SLAS120=103M

LIST=445 OBSERVATIONS PER TRACE=423 GENERATED BY DATA

ALPHABETIC LIST OF VARIABLES

#	VARIABLE	TYPE	LENGTH	POSITION	FORMAT	INFORMAT	LABEL
2	CNUM	CHAR	6	32			
1	DEC	NUM	8	4			
3	HEX	CHAR	8	16			
4	ROMAN	CHAR	6	24			
5	WORDS	CHAR	15	30			

SOURCE STATEMENTS

```

DATA EXAMPLE;
  INPUT DEC 80;
  CNUM = PUT(DEC,4.);
  HEX = PUT(DEC,HEX.);
  ROMAN = PUT(DEC,ROMAN.);
  WORDS = PUT(DEC,WORDS);
CARDS;
  
```


2. CREATING FORMATS

(1) The FORMAT Procedure

The FORMAT procedure is used to create user-defined formats.

These user-defined formats can be used:

- o in a PUT statement
- o in a FORMAT statement with a procedure.

(2) PROC FORMAT Options;

Selected options:

PRINT

DDNAME=ddname

Statements used with PROC FORMAT:

VALUE name (options)

range 1 = 'label1'

range 2 = 'label2'

... ;

PICTURE name (options)

range 1 = 'picture1' (options)

range 2 = 'picture2' (options)

(3) The VALUE Statement

Assigning values:

- o single numbers

VALUE Q 1 = 'AGREE' 2 = 'DISAGREE'

- o ranges of numbers

VALUE AGEFMT 0-12 = CHILD
13-19 = 'TEEN'
20-HIGH = 'ADULT';

- o several values

VALUE SEXFMT 1 = 'FEMALE'
2 = 'MALE'
0, 3-9 = 'MISCODED'

- o character values and ranges of characters

VALUE \$GRADE A = 'GOOD'
B-D = 'FAIR'
E = 'POOR'
I,U = 'SEE INSTRUCTOR';

- o character values with special characters.

VALUE \$CODEX 'A*1' = 'FIRST'
'A*2' = 'SECOND'
A__ = 'MORE THAN TWO';

Assigning labels:

- o maximum length of 40 characters
- o labels must be enclosed in single quotes.

(4) Date, Time, and Datetime Informat and Formats

Name of format or informat	Example	Informat, format, or both
DATEw.	04JUL1976	both
YYMMDDw.	76-07-04	both
MMDDYYw.	7/4/76	both
DDMMYYw.	4/7/76	both
MONYYw.	Jul76	both
YYQw.	76Q3	both
WEEKDATEw.	Monday, July 4, 1976	format
WORDDATEw.	July 4, 1976	format
HHMMw.d	23:45	format
HOURw.d	23	format
MMSSw.d	45:23.4	format
MSECw.	TIME MIC values	informat
PDTIMEw.	packed-decimal time from RMF records	informat
RMFDURw.	RMF time interval measurements	informat
TIMEw.d	23:45:23.5	both
TODw.	23:45:23.4	format
TUw.	timer units	informat
DATETIMEw.	OJUL1976:23:45:23.5	both
RMFSTAMPw.	RMF time-date field	informat
SMFSTAMPw.	SMF time-date field	informat
TODSTAMPw.	8-byte time-of-day stamp	informat

Note: See pages 409-421 in the SAS User's Guide:
Basics, 1982 Edition.

(5) Date, Time, and Datetime Values

Example : A company wants to determine the length of
Employment in years for each of its
employees who resigned in 1981.

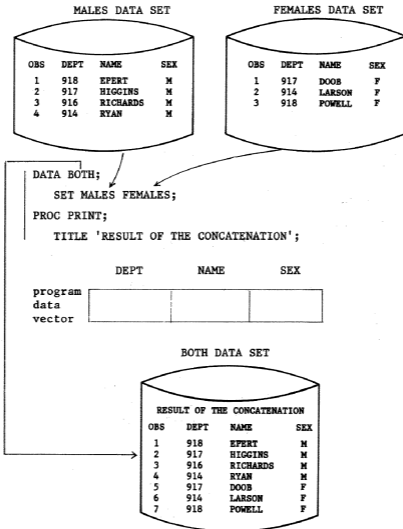
```
//SPSSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
DATA RESIGNED;
  INPUT NAME $ 10. +1 FIRSTDAY MMDDYY8.
        +1 LASTDAY MMDDYY8.;
  DAYS=LASTDAY-FIRSTDAY;
CARDS;
ARONTH, J. 12/01/73 4/30/81
DESTER, L. 7/14/61 12/31/81
HARLIN, M. 8/03/77 6/15/81
;
PROC PRINT;
/*
//
```

- o Remember that SAS date and time values have implicit units.

OBS	NAME	SAS		
		FIRSTDAY	LASTDAY	DAYS
1	ARONTH, J.	5083	7790	2707
2	DESTER, L.	300	8035	7475
3	HARLIN, M.	6624	7636	1412

(6) Concatenating SAS Data Sets

Example : Combine the MALES and FEMALES data sets
into one data set.



(7) Interleaving SAS Data Sets

Example : Combine the MALES and FEMALES data sets such that the resulting data set has its observations arranged in alphabetical order.

```
PROC SORT DATA=MALES; BY NAME;
PROC SORT DATA=FEMALES; BY NAME;
```

MALES DATA SET

OBS	DEPT	NAME	SEX
1	918	EPERT	M
2	917	HIGGINS	M
3	916	RICHARDS	M
4	914	RYAN	M

FEMALES DATA SET

OBS	DEPT	NAME	SEX
1	917	DOOB	F
2	914	LARSON	F
3	918	POWELL	F

```
DATA BOTHSORT;
  SET MALES FEMALES;
  BY NAME;
PROC PRINT;
  TITLE 'RESULT OF INTERLEAVING';
```

	DEPT	NAME	SEX
program			
data			
vector			

BOTHSORT DATA SET

RESULT OF INTERLEAVING			
OBS	DEPT	NAME	SEX
1	917	DOOB	F
2	918	EPERT	M
3	917	HIGGINS	M
4	914	LARSON	F
5	918	POWELL	F
6	916	RICHARDS	M
7	914	RYAN	M

(8) Match Merging

Data sets with unequal numbers of observations

Example : Suppose Larson is missing from the SALARY data set. (Note: Both data sets are already sorted by name.)

GENERAL DATA SET

OBS	DEPT	NAME	SEX
1	917	DOOB	F
2	918	EPERT	M
3	917	HIGGINS	M
4	914	LARSON	F
5	918	POWELL	F
6	916	RICHARDS	M
7	914	RYAN	M

SALARY DATA SET

OBS	NAME	NETPAY	GROSSPAY
1	DOOB	169.06	272.29
2	EPERT	224.36	310.40
3	HIGGINS	777.50	1235.46
4	POWELL	189.39	271.54
5	RICHARDS	219.27	352.84
6	RYAN	291.56	399.20

```
DATA MERGED;  
MERGE GENERAL SALARY;  
BY NAME;  
PROC PRINT;  
TITLE 'MATCH MERGING';  
TITLE2 'UNEQUAL NUMBERS OF OBSERVATIONS';
```

	DEPT	NAME	SEX	NETPAY	GROSSPAY
program data vector					

MERGED DATA SET

MATCH MERGING UNEQUAL NUMBERS OF OBSERVATIONS					
OBS	DEPT	NAME	SEX	NETPAY	GROSSPAY
1	917	DOOB	F	169.06	272.29
2	918	EPERT	M	224.36	310.40
3	917	HIGGINS	M	777.50	1235.46
4	914	LARSON	F	.	.
5	918	POWELL	F	189.39	271.54
6	916	RICHARDS	M	219.27	352.84
7	914	RYAN	M	291.56	399.20

(9) Match Merging

Identical variable names

CLOTHES DATA SET

OBS	DATE	SALES
1	18OCT82	223.93
2	19OCT82	387.82
3	20OCT82	229.28
4	21OCT82	318.32
5	22OCT82	519.07

EQUIP DATA SET

OBS	DATE	SALES
1	18OCT82	492.28
2	19OCT82	228.20
3	20OCT82	542.98
4	21OCT82	325.02
5	22OCT82	733.60

```
DATA ALLSALES;  
MERGE CLOTHES EQUIP;  
BY DATE;  
  
PROC PRINT;  
TITLE 'MERGING DATA SETS';  
TITLE2 'WITH IDENTICAL VARIABLE NAMES';
```

```
DATE          SALES  
program  
data  
vector
```

--	--

ALLSALES DATA SET

MERGING DATA SETS WITH IDENTICAL VARIABLE NAMES		
OBS	DATE	SALES
1	18OCT82	492.28
2	19OCT82	228.20
3	20OCT82	542.98
4	21OCT82	325.02
5	22OCT82	733.60

(10) Match Merging

Identical variable names

Example : Repeat the preceding example using the RENAME data set option.

CLOTHES DATA SET

OBS	DATE	SALES
1	18OCT82	223.93
2	19OCT82	387.82
3	20OCT82	229.28
4	21OCT82	318.32
5	22OCT82	519.07

EQUIP DATA SET

OBS	DATE	SALES
1	18OCT82	492.28
2	19OCT82	228.20
3	20OCT82	542.98
4	21OCT82	325.02
5	22OCT82	733.60

```
DATA ALLSALES;  
MERGE CLOTHES (RENAME (SALES=CL_SALES))  
      EQUIP (RENAME=(SALES=EQ_SALES));  
BY DATE;  
PROC PRINT;  
TITLE 'RESULT OF MERGING DATA SETS';  
TITLE2 'WITH IDENTICAL VARIABLE NAMES';  
TITLE3 'USING THE RENAME DATA SET OPTION';
```

	DATE	CL_SALES	EQ_SALES
--	------	----------	----------

program
data
vector

--	--	--	--

ALLSALES DATA SET

RESULT OF MERGING DATA SETS WITH IDENTICAL VARIABLE NAMES USING THE RENAME DATA SET OPTION			
OBS	DATE	CL_SALES	EQ_SALES
1	18OCT82	223.93	492.28
2	19OCT82	387.82	228.20
3	20OCT82	229.28	542.98
4	21OCT82	318.32	325.02
5	22OCT82	519.07	733.60

(11) UPDATE Application

Example : Add new employees and make changes to existing values for old employees in the PAYROLL data set.

(Note: Both data sets have already been sorted by Name.)

PAYROLL DATA SET

OBS	DEPT	NAME	SEX	NETPAY	GROSSPAY
1	917	DOOB	F	169.06	272.29
2	918	EPERT	M	224.36	310.40
3	917	HIGGINS	M	777.50	1235.46
4	914	LARSON	F	215.47	283.92
5	918	POWELL	F	189.39	271.54
6	916	RICHARDS	M	219.27	352.84
7	914	RYAN	M	291.56	399.20

DATA NEWINFO;

INPUT DEPT 1-3 NAME \$ 5-12 SEX \$14
NETPAY 16-22 GROSSPAY 24-30;

CARDS;

POWELL 221.75 310.62
916 SERPANT M 207.22 398.65
DOOB 191.65 252.57
918 ARCHER F 315.17 420.00

PROC SORT DATA=NEWINFO;
BY NAME;

PROC PRINT DATA=NEWINFO;

NEWINFO DATA SET

OBS	DEPT	NAME	SEX	NETPAY	GROSSPAY
1	918	ARCHER	F	315.17	420.00
2	.	DOOB		191.65	252.57
3	.	POWELL		221.75	310.62
4	916	SERPANT	M	207.22	398.65

(12) UPDATE Application

Example : Update the PAYROLL data set with the
NEWINFO data set.

```
DATA PAYROLL2;  
  UPDATE PAYROLL NEWINFO;  
  BY NAME;  
PROC PRINT;  
  TITLE 'PAYROLL2 DATA SET';
```

PAYROLL2 DATA SET

OBS	DEPT	NAME	SEX	NETPAY	GROSSPAY
1	918	ARCHER	F	315.17	420.00
2	917	DOOB	F	191.65	252.57
3	918	EPERT	M	224.36	310.40
4	917	HIGGINS	M	777.50	1235.46
5	914	LARSON	F	215.47	283.92
6	918	POWELL	F	221.75	310.62
7	916	RICHARDS	M	219.27	352.84
8	914	RYAN	M	291.56	399.20
9	916	SERPANT	M	207.22	398.65



(13) SAS Array

Example: Create the 4 Variables with a SAS Array.

```
//SPSSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
DATA CONVERT;
INPUT N1-N5 ;
ARRAY PERCENT N1-N5;
TOTAL=SUM(OF N1-N5);

DO OVER PERCENT;
PERCENT=ROUND((PERCENT/TOTAL*100),1);
END;
TOTAL=SUM(OF N1-N5);
CARDS;
11 5 26 9 49
19 5 4 36 37
43 3 9 8 36
6 25 28 34 7
44 13 6 5 32
;
PROC PRINT;
TITLE ARRAY SAMPLE;
/*
//
```

10:59 FRIDAY, NOVEMBER 16, 1

ARRAY SAMPLE						
OPS	N1	N2	N3	N4	N5	TOTAL
1	11	5	26	9	49	100
2	19	5	4	36	37	101
3	43	3	9	8	36	99
4	6	25	28	34	7	100
5	44	13	6	5	32	100

CHAPTER 3. 기 초 통 계 처 리

CHAPTER 3. 기초통계처리

1. CHART

2. PLOT

3. Discriptive statistic

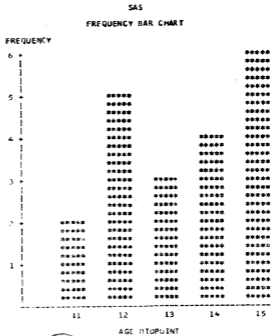
- 1) CORR computes bivariate correlations and other measures of association for continuous variables.
- 2) FREQ print tables of frequency counts, cross tabulations, and bivariate measures of association for categorical variables.
- 3) MEANS computes and prints means and other descriptive statistics.
- 4) UNIVARIATE computes univariate statistics including quantiles.

* Categorical data

1. CHART Procedure

(1) 기본형

```
PROC CHART;  
VBAR AGE;
```



(2) 일반형

```
PROC CHART options;  
    BY          variables;  
    VBAR        variables/options;  
    HBAR        variables/options;  
    BLOCK       variables/options;  
    PIE         variables/options;  
    STAR        variables/options;
```

(3) Options

1) Determining the values represented

TYPE = code (FREQ CFREQ PCT CPCT SUM or MEAN)

SUMVAR = variable

FREQ = variable

2) Grouping among and within bars

GROUP = variable

SUBGROUP = variable name

3) Classifying the observations into bars

DISCRETE MIDPOINTS LEVELS

4) Formatting the chart

NOSPACE SYMBOL = 'character' MISSING AXIS=values

NOSTAT FREQ CFREQ PERCENT CPERCENT SUM MEAN

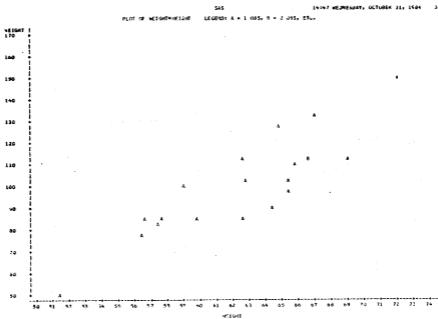
2. PLOT PROCEDURE

(1) 기본형

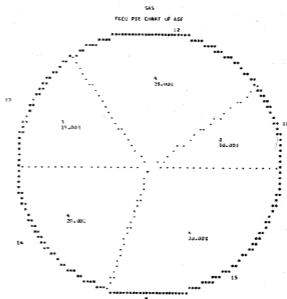
```
PROC PLOT;
```

```
    PLOT HEIGHT * WEIGHT;
```

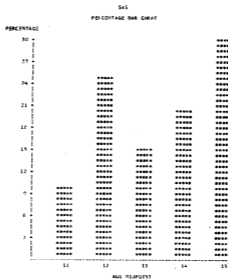
1) PLOTS



2) PROC CHART;
 PIE AGE;



3) PROC CHART;
 VBAR AGE;



(2) 일반형

PROC PLOT DATA=dataset UNIFORM NOLEGEND;

BY variables;

PLOT requests/options;

- o request : vertical *horizontal=('character') or
variable
- o scale of axis options
VAXIS HAXIS VZERO VREVERSE
- o plot size options
VPOS=n HPOS=n VSPACE=n HSPACE=n
- o overlaying plots
OVERLAY
- o contour plots
CONTOUR = n (n ranges 1 to n).

3. CORR procedure

(1) Product-moment Correlation (Pearson)

1) True product moment correlation

$$\rho_{xy} = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x) \text{var}(y)}} \\ = \frac{E(x-E(x))(y-E(y))}{\sqrt{E(x-E(x))^2 E(y-E(y))^2}}$$

2) Sample correlation estimates

$$= \frac{(x-\bar{x})(y-\bar{y})}{\sqrt{(x-\bar{x})^2 (y-\bar{y})^2}}$$

- 3) Spearman's rank order correlation coefficient;
 nonparametric measure that calculates the correlation of the ranks of the data

$$\rho = \frac{\sum (r_i - \bar{r})(s_i - \bar{s})}{\sqrt{\sum (r_i - \bar{r})^2 \sum (s_i - \bar{s})^2}}$$

- 4) Kendall's tau- τ
 ; measure calculated from concordances and discordances

$$\tau = \frac{\sum_i \sum_j \text{Sgn}(x_i - x_j) \text{Sgn}(y_i - y_j)}{\sqrt{((n-1)/2 - t_i(t_i-1)) (n-1) - \sum_i (u_i(u_i-1))}}$$

(2) 기본형

```
PROC CORR;
  VAR HEIGHT WEIGHT;
  WITH AGE;
```

VARIABLE	N	MEAN	STD DEV	SUM	MINIMUM	MAXIMUM
HEIGHT	20	175.5000000	5.24128439	3501.0000000	81.0000000	172.0000000
WEIGHT	20	130.9500000	22.78136996	2711.0000000	50.5000000	155.0000000
AGE	20	15.4000000	1.50016766	308.0000000	11.0000000	18.0000000

CORRELATION COEFFICIENTS / P&B > (R) UNDER NO&NO=0 / N = 20

	HEIGHT	WEIGHT
AGE	0.11600	0.73000
	0.0001	0.0002

(3) 일반형

```
PROC CORR options;  
    VAR      variables;  
    WITH     variables;  
    WEIGHT   variables;  
    FREQ     variables;  
    BY       variables;
```

(4) Options

```
SPEARMAN   BEST      NOMISS  
KENDALL    NOSIMPLE  SSCP  
PEARSON    NOPRINT   COV  
RANK       NOPROB    NCORR
```

OUTP= SAS dataset: requests that CORR create a new SAS
data set containing Pearson correlation

OUTS= SAS dataset: requests that CORR create a new SAS
data set containing Spearman correlations

OUTK= SAS dataset: requests that CORR create a new SAS
data set containing Kendall correlations.

4. FREQ procedure. (discrete variable에 사용)

(1) 기본형

PROC FREQ

TABLE AGE *SEX;

```
//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
```

```
DATA SAMPLE;
INPUT NAME $ 1-10 SEX $ 11 AGE 13-14
HEIGHT 16-19 WEIGHT 21-25;
```

```
CARDS;
KYUNGEUI M 15 65.7 110.9
JOHN M 12 59.0 99.5
JAMES M 12 57.3 83.0
ALFERD M 14 69.0 112.5
WILLIAM M 15 66.5 112.0
JEFFREY M 13 62.5 84.0
RONALD M 15 67.0 133.0
THOMAS M 11 57.5 85.5
PHILIP M 16 72.0 150.0
ROBELT M 12 64.8 128.0
HENRY M 14 65.3 102.5
JANET F 15 62.5 112.5
JOYCE F 11 51.3 50.5
JUDY F 14 64.3 90.0
CAROL F 14 62.8 102.5
JANE F 12 59.0 84.5
LOUISE F 12 56.3 77.0
BARBARA F 13 65.3 98.0
MARY F 15 66.5 112.0
ALICE F 13 56.5 84.0
```

```
PROC FREQ;
TABLES AGE*SEX;
```

```
/*
//
```

SAS
TABLE OF AGE BY SEX

AGE	SEX		TOTAL
	F	M	
11	1	1	2
	5.00	5.00	10.00
	50.00	50.00	
	11.00	9.00	
12	2	3	5
	10.00	15.00	25.00
	40.00	60.00	
	22.22	27.27	
13	2	1	3
	10.00	5.00	15.00
	40.00	20.00	
	22.22	9.09	
14	2	2	4
	10.00	10.00	20.00
	50.00	50.00	
	22.22	10.10	
15	2	3	5
	10.00	15.00	25.00
	40.00	60.00	
	22.22	27.27	
16	0	1	1
	0.00	5.00	5.00
	0.00	100.00	
	0.00	9.09	
TOTAL	9	11	20
	45.00	55.00	100.00



(2) 일반형

PROC FREQ options;
TABLES requests/options;
WEIGHT variables;
BY variables;

(3) TABLES Options

EXPECTED requests that the expected cell frequency under the hypothesis of independence (or homogeneity) be printed.

DEVIATION requests that FREQ print deviation of the cell frequency from the expected value.

CELLCHI2 requests that FREQ print the cell's contribution to the total χ^2 statistic. This is computed as (frequency-expected)**2/expected and is approximately distributed χ^2 with 1 df.

CHISQ requests a chi-square(χ^2) test of homogeneity or independence for each two-way table requested in a TABLES statement. For 2 by 2 tables, Fisher's Exact Test is performed. The formula for χ^2 is given in Chapter17,

"Introduction to Descriptive Statistics."

ALL requests the basic set of measures of association popularized by Goodman and Kruskal for two-way tables, including some of the standard errors, the contingency coefficient, Cramer's V, gamma, Kendall's tau- ζ , Stuart's tau-c, Somer's D, and lambda asymmetric and Spearman correlations. Of course, not all statistics are appropriate for the data in a given table, (See ALL Option: Measures of Association below.)

NOFREQ suppresses printing the cell frequencies for a crosstabulation.

NOPERCENT suppresses printing the cell percentages for a crosstabulation.

NOROW suppresses printing the row percentages in cells of a crosstabulation.

NOCOL suppresses printing the column percentages in cells of a crosstabulation.

CUMCOL requests cumulative column percentages be printed in the cells.

LIST prints two-way to n-way tables in a list format rather than as crosstabulation tables. Expected cell frequencies are not printed when LIST is specified, even if EXPECTED is specified.

NOCUM suppresses the cumulative frequencies, percentages, and cumulative percentage columns for one-way frequencies and frequencies in list format when the LIST option is included.

MISSING asks FREQ to consider missing values like other values in calculations of percents and other statistics.

SPARSE causes the procedure to write out or print information about all possible combinations of levels of the variables in the table request, even when some com-

binations of levels do not occur in the data.

This option affects printouts under the LIST option and output data sets.

NOPRINT suppresses all printed output except that controlled by CHISQ and ALL.

5. MEANS procedure.

(1) 기본형

```
PROC MEANS;  
    VAR HEIGHT WEIGHT;  
    BY SEX;
```

VARIABLE	N	MEAN	STANDARD DEVIATION	SAS		STD ERROR OF MEAN	SUM	VARIANCE	C.V.
				MINIMUM VALUE	MAXIMUM VALUE				
HEIGHT	9	60.70918889	5.51882792	51.30000000	66.50000000	1.67277566	565.30000000	29.18261111	0.283
WEIGHT	9	70.31111111	15.28391372	50.30000000	112.50000000	6.44230457	631.00000000	179.79111111	21.911

HEIGHT	11	64.23636364	4.72468576	57.00000000	72.00000000	1.42466666	706.00000000	22.32454545	7.233
WEIGHT	11	109.17272727	21.91344189	83.00000000	150.00000000	6.48456679	1206.00000000	437.82181818	16.704

(2) 일반형

```
PROC MEANS options;  
    BY variables;  
    VAR variables;  
    FREQ variables;  
    ID variables;  
    OUTPUT OUT = SAS dataset keyword = names...;
```

(3) PROC MEANS options;

N	the number of observations on which calculations are based
NMISS	the number of missing values
MEAN	the mean
STD	the standard deviation
MIN	the smallest value
MAX	the largest value
RANGE	the range
SUM	the sum
VAR	the variance
USS	the uncorrected sum of squares
CSS	the corrected sum of squares
STDERR	the standard error of the mean
CV	the coefficient of variation (percent)
SKEWNESS	the measure of skewness
KURTOSIS	the measure of kurtosis
T	the student's t value for testing the hypothesis that the population mean is 0
PRT	the probability of a greater absolute value of Student's t.

(4) STATISTICS

$$\text{VAR} = \sum (x-\bar{x})/n-1$$

$$\text{SKEWNESS} = m_3/m_2^{2/2}$$

$$\text{KURTOSIS} = m_4/m_2^2$$

$$T = \bar{x}/\sqrt{\text{VAR}/N}$$

$$\text{PRT} = \text{Pro}(|t| > T)$$

6. UNIVARIATE procedure

(1) 기본형

```
PROC UNIVARIATE;  
VAR HEIGHT AGE;  
BY SEX;
```

UNIVARIATE procedure에서는 다른 기초통계 procedure에서
계산되지 않는 백분위수가 계산된다.

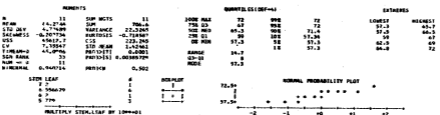
```
//SPCSAS30 JOB CLASS=V  
//SAS EXEC SAS  
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR  
//SYSIN DD *  
DATA SAMPLE;  
INPUT NAME $ 1-10 SEX $ 11 AGE 13-14  
HEIGHT 16-19 WEIGHT 21-25;  
CARDS;  
KYUNGEUI M 15 65.7 110.9  
JOHN M 12 59.0 99.5  
JAMES M 12 57.3 83.0  
ALFERD M 14 69.0 112.5  
WILLIAM M 15 66.5 112.0  
JEFFREY M 13 62.5 84.0  
RONALD M 15 67.0 133.0  
THOMAS M 11 57.5 85.5  
PHILIP M 16 72.0 150.0  
ROBELT M 12 64.8 128.0  
HENRY M 14 65.3 102.5  
JANET F 15 62.5 112.5  
JDYCE F 11 51.3 50.5  
JUDY F 14 64.3 90.0  
CAROL F 14 62.8 102.5  
JANE F 12 59.8 84.5  
LOUISE F 12 56.3 77.0  
BARBARA F 13 65.3 98.0  
MARY F 15 66.5 112.0  
ALICE F 13 56.5 84.0  
;  
PROC UNIVARIATE FREQ PLOT NORMAL,  
VAR HEIGHT AGE;  
/*  
//
```

SAS

14147 WEDNESDAY, OCTOBER 21, 1964

UNWEIGHTED

VARIABLE=HEIGHT



FREQUENCY TABLE

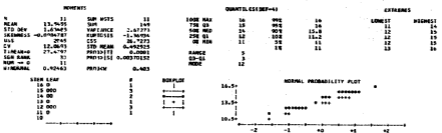
PERCENTS			PERCENTS			PERCENTS		
VALUE	COUNT	CUM	VALUE	COUNT	CUM	VALUE	COUNT	CUM
167.5	1	9.1	172.0	1	18.2	176.5	1	27.3
172.0	2	18.2	176.5	1	27.3	180.5	1	36.4
176.5	1	27.3	180.5	1	36.4			

SAS

14147 WEDNESDAY, OCTOBER 21, 1964

UNWEIGHTED

VARIABLE=AGE



FREQUENCY TABLE

PERCENTS			PERCENTS			PERCENTS		
VALUE	COUNT	CUM	VALUE	COUNT	CUM	VALUE	COUNT	CUM
10	5	45.5	10	1	9.1	10	1	9.1
11	1	54.5	11	1	18.2	11	1	27.3
12	3	77.3	12	2	36.4			

(2) STATISTICS

N	the number of observations on which the calculations were based
NMISS	the number of missing values
NOBS	the number of observations
MEAN	the mean
SUM	the sum
STD	the standard deviation
VAR	the variance
SKEWNESS	skewness
KURTOSIS	kurtosis
SUMWGT	the sum of the weights
MAX	the largest value
MIN	the smallest value
RANGE	the range
Q3	the upper quartile or the seventy-fifth percentile
MEDIAN	the median or the fiftieth percentile
Q1	the lower quartile or the twenty-fifth percentile
QRANGE	the difference between the upper and lower quartiles, that is, $Q3-Q1$
P1	the first percentile
P5	the fifth percentile
P10	the tenth percentile
P90	the ninetieth percentile
P95	the ninety-fifth percentile

P99 the ninety-ninth percentile
MODE the most frequent value. If the mode is not
 unique, the smallest mode is used

(3) 일반형

PROC UNIVARIATE options;
 VAR variables;
 BY variables;
 FREQ variables;
 WEIGHT variables;
 ID variables;
 OUTPUT OUT = SAS data set key word = names,...;

(4) PROC UNIVARIATE options;

NOPRINT suppresses all printed output. NOPRINT can
 be used when the only purpose for executing
 the procedure is to create new data sets.
PLOT causes UNIVARIATE to produce a stem-and-leaf
 plot (or a horizontal bar char), a box plot,
 and a normal probability plot.
FREQ requests a frequency table consisting of the
 variable values, frequencies, percentages,
 and cumulative percentages.

NORMAL causes UNIVARIATE to compute a test statistic for the hypothesis that the input data come from a normal distribution. The probability of a more extreme value of the test statistic is also printed.

DEF=value specifies which of the four definitions given below in the section Computational Method is to be used to calculate percentiles. The DEF value may be 1,2,3,4, or 5, If DEF= is omitted, definition 4 is used.

CHAPTER 4. 회귀분석과 분산분석에 관련된 기법

CHAPTER 4. 회귀분석과 분산분석에 관련된 기법

A. 회귀분석 (Regression Analysis)

1. 회귀분석

- (1) 회귀분석은 종속변수 (dependent variable) 와 독립변수 (independent variables) 들간의 관계를 알아보기 위해 사용된다. 독립변수는 설명변수 (explanatory variable) 라고도 한다. SAS의 REG procedure 를 이용하면 선형 회귀분석을 할 수 있다.

- (2) 선형회귀모형 : $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$

2. REG procedure

- (1) PROC REG options;

MODEL dependents = independents /options;

WEIGHT variable;

ID variable;

OUTPUT OUT = SAS dataset keyword = names...;

TEST linear equations...;

- (2) PROC REG Statement

대표적인 예로는

PROC REG DATA = A; SAS dataset A를 이용한다.

PROC REG; 가장 최근에 만들어진 SAS dataset를 이용

표준화된 printout을 만든다.

PROC REG ALL; 가장 자세한 printout을 만든다.

- (3) MODEL statement의 대표적인 예로는

```
MODEL Y=X1 X2 /XPX I ;
```

선형모형 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$ 와 관련된 각종 통계계산을 하고 option으로 $X'X$ 와 $(X'X)^{-1}$ 를 printout 한다.

```
MODEL Y=X1 X2 /COLLIN ;
```

multicollinearity 분석을 option으로 수행한다.

```
MODEL Y=X1 X2 /PARTIAL
```

종속변수와 독립변수의 partial scatter plot 을 만든다.

- (4) WEIGHT variable; 가장 최소자승법에 의한 회귀분석을 한다.
- (5) ID variable; 각 관측치를 표시하는 변수를 이용할 때 사용한다.
- (6) OUTPUT statement의 대표적인 예로는

```
OUTPUT OUT = B
```

```
PREDICTED = YHAT
```

```
RESIDUAL = YRESID ;
```

종속변수의 predicted values 을 YHAT, 잔여분 (residual) 을

YRESID로 이름짓고 SAS dataset B에 넣는다.

- (7) TEST statement의 대표적인 예로는

```
MODEL Y = X1 X2 X3 X4 ;
```

선형모형 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \epsilon$ 에 관련된 회귀

분석에서 통계적가설 $\beta_3 = \beta_4 = 0$ 을 검정한다.

- ※ 그 밖의 자세한 건반적인 소개와 설명을 SAS User's Guide ;
Statistics pp. 39 - 83 을 참조할 것.

(8) Sample

다음은 83년도 도시가계조사 대상중 15가구에 대한 X (소득),

Y 1 (소비지출), Y 2 (교육교양오락비)의 자료이다.

```
//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SYSIN DD *
DATA HOUSE;
INPUT X Y1 Y2;
CARDS;
66618 118319 8534
122469 136736 9432
171358 157592 10050
219049 184560 13376
270043 208717 15605
319007 243169 24881
370755 271723 26245
420115 305076 34969
472620 339056 38686
519782 355504 34544
572612 383370 42598
640551 423327 53675
741887 461573 54641
844780 520846 53965
1221396 666157 91260
;
PROC REG ;
MODEL Y1=X;
PROC REG ;
MODEL Y2=X;
/*
//
```

SAS

DEP VARIABLE: Y1

SOURCE	DF	SUM OF SQUARES	MEAN SQUARE	F VALUE	PROB>F
MODEL	1	333741990533	333741990533	1692.873	0.0001
ERROR	13	2562888776	197145290		
C TOTAL	14	336304879309			
ROOT MSE		14040.844	R-SQUARE	0.9924	
DEP MEAN		318382	ADJ R-SQ	0.9918	
C.V.		4.410067			

VARIABLE	DF	PARAMETER ESTIMATE	STANDARD ERROR	T FOR HO: PARAMETER=0	PROB > T
INTERCEP	1	85641.266	6718.686	12.747	0.0001
X	1	0.500658	0.012168	41.145	0.0001

SAS

DEP VARIABLE: Y2

SOURCE	DF	SUM OF SQUARES	MEAN SQUARE	F VALUE	PROB>F
MODEL	1	7114011490	7114011490	449.167	0.0001
ERROR	13	205897037	15838234		
C TOTAL	14	7319908527			
ROOT MSE		3979.728	R-SQUARE	0.9719	
DEP MEAN		34164.067	ADJ R-SQ	0.9697	
C.V.		11.64887			

VARIABLE	DF	PARAMETER ESTIMATE	STANDARD ERROR	T FOR HO: PARAMETER=0	PROB > T
INTERCEP	1	184.047	1904.340	0.097	0.9245
X	1	0.073096	0.003448964	21.194	0.0001

SAS

DEP VARIABLE: OXY

SOURCE	DF	SUM OF SQUARES	MEAN SQUARE	F VALUE	PROB>F
MODEL	6	202.659	33.776418	4.860	0.0175
ERROR	9	62.55317	6.950435		
C TOTAL	15	265.212			
ROOT MSE		2.636368		0.7641	
DEP MEAN		46.284937	R-SQUARE	0.6069	
C.V.		5.659552	ADJ R-SQ		

VARIABLE	DF	PARAMETER ESTIMATE	STANDARD ERROR	T FOR H0: PARAMETER=0	>PROB > T
INTERCEPT	1	49.360100	26.377029	1.871	0.0941
RUNTIME	1	-0.748050	0.443428	-1.687	0.1259
AGE	1	-0.218662	0.142014	-1.540	0.1580
WEIGHT	1	0.103334	0.098252	1.041	0.3250
RUNPULSE	1	-0.314025	0.279587	-1.123	0.2904
MAXPULSE	1	0.454630	0.343396	1.324	0.2182
RSTPULSE	1	-0.331326	0.178302	-1.848	0.0977

COLLINEARITY DIAGNOSTICS VARIANCE PROPORTIONS

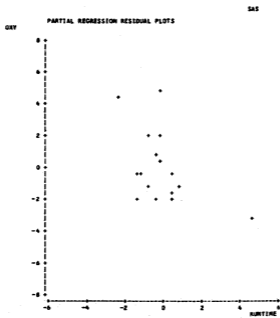
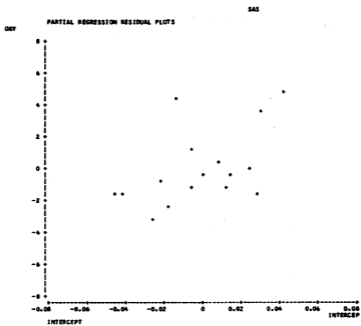
NUMBER	EIGENVALUE	CONDITION INDEX	PORTION INTERCEPT	PORTION RUNTIME	PORTION AGE	PORTION WEIGHT	PORTION MAXPULSE	PORTION RSTPULSE
1	6.939	1.000	0.0000	0.0004	0.0002	0.0001	0.0000	0.0001
2	0.031822	14.767	0.0008	0.3719	0.0299	0.0183	0.0031	0.0001
3	0.015039	21.480	0.0000	0.0626	0.3548	0.1035	0.0030	0.0001
4	0.009015	26.455	0.0001	0.2055	0.0001	0.3260	0.0013	0.0006
5	0.003643	44.895	0.0227	0.0894	0.4816	0.1542	0.0097	0.0040
6	0.0005113	116.497	0.8201	0.1042	0.1273	0.3458	0.0818	0.0105
7	0.00074964	304.248	0.1564	0.1661	0.0362	0.0520	0.0970	0.9847

(9) Sample

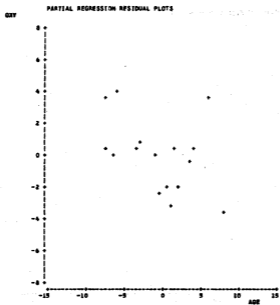
Aerobic Fitness Prediction :

Aerobic fitness 는 산소 소비능력으로 측정되는데 이 측정 방법은 비용이 많이 들고 불편하므로 그 대신 간단한 운동측정으로 산소 소비능력을 예측하고자 한다. 회귀분석에 사용된 종속변수는 OXY (oxgen uptake rate 산소소비율), 설명변수로는 RUNTIME (1.5 마일 주행시간), AGE (나이), WEIGHT (체중), RUMPULSE (주행중의 맥박수), MAXPULSE (주행중의 최대맥박수), RSTPULSE (휴식중의 맥박수) 포함되었다.

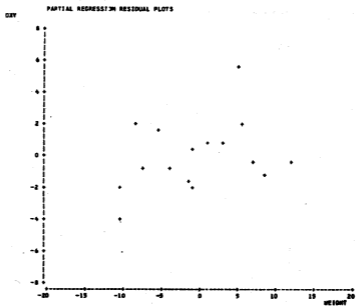
```
//SPCSAS30 JOB CLASS=M
//SAS EXEC SAS
//SYSIN DD *
DATA FITNESS;
INPUT AGE WEIGHT OXY RUNTIME RSTPULSE
      RUMPULSE MAXPULSE;
CARDS;
44 89.47 44.609 11.37 62 170 182 40 75.07 45.313 10.07 62 185 185
44 85.84 54.297 8.65 45 156 168 42 68.15 59.574 8.17 40 166 172
38 89.02 49.874 9.22 55 178 180 47 77.45 44.811 11.63 58 176 176
40 75.98 45.681 11.95 70 176 180 43 81.19 49.094 10.85 64 162 170
44 81.42 39.442 18.08 63 174 176 38 84.87 60.055 8.63 48 170 186
44 73.03 50.541 10.13 45 168 168 45 87.66 37.388 14.03 56 186 192
45 66.45 44.754 11.12 51 176 176 47 79.15 47.273 10.60 47 162 164
54 83.12 51.855 10.33 50 166 170 49 81.42 49.156 8.95 44 180 185
51 69.63 46.836 10.95 57 168 172 51 77.91 46.672 10.00 48 162 168
48 91.63 46.774 10.25 48 162 164 49 73.37 50.388 10.08 67 168 168
57 73.37 39.407 12.63 58 174 176 54 79.38 46.080 11.17 62 156 165
52 76.32 45.441 9.63 48 164 166 50 70.87 54.625 8.92 48 146 155
51 67.25 45.118 11.08 48 172 172 54 91.63 39.203 12.88 44 168 172
51 73.71 45.790 10.47 59 186 188 57 59.08 50.545 9.93 49 148 155
49 76.32 48.673 9.40 56 186 188 48 61.24 47.920 11.50 52 170 176
52 82.78 47.467 10.50 53 170 172
;
PROC REG OUTEST=EST;
MODEL OXY=RUNTIME AGE WEIGHT RUMPULSE
      MAXPULSE RSTPULSE/PARTIAL COLLIN;
/*
//
```



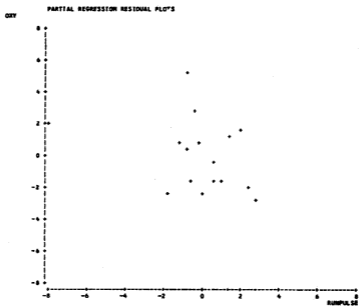
SAS



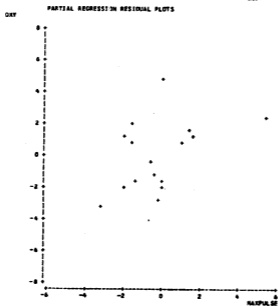
SAS



SAS

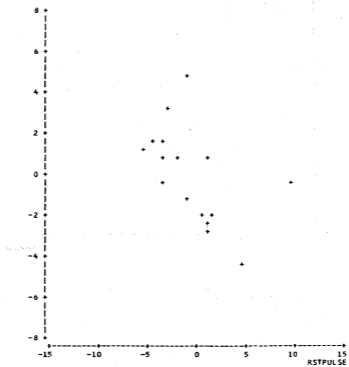


SAS



DXY

PARTIAL REGRESSION RESIDUAL PLOTS



3. STEPWISE procedure

(1) PROC STEPWISE option;

MODEL dependents = independent/options;
WEIGHT variable;
BY variable;

(2) PROC STEPWISE option

NOINT	prevents the procedure from automatically including an intercept term in the model.
FORWARD F	requests the forward-selection technique.
BACKWARD B	requests the backward-elimination technique.
STEPWISE	requests the stepwise technique, the default.
MAXR	requests the maximum R^2 improvement technique.
MINR	requests the minimum R^2 improvement technique.
SLENTRY = value SLE = value	specifies the significance level for entry into the model used in the forward-selection and stepwise techniques. If SLENTRY = is omitted, STEPWISE uses the SLENTRY = value .50 for forward selection, .15 for stepwise.

SLSTAY = value specifies the significance level for
SLS = value staying in the model for the backward
elimination and stepwise techniques.
If it is omitted, STEPWISE uses the
SLSTAY = value .10 for backward
elimination, .15 for stepwise.

INCLUDE = n forces the first n independent
variables always to be included in the
model. The selection techniques are
performed on the other variables in
the MODEL statement.

START = s is used to begin the comparing-and-
switching process for a model contain-
ing the first s independent variables
in the MODEL statement, where s is the
START value. Consequently, no model
is evaluated that contains fewer than
s variables. This applies only to the
MAXR or MINR methods.

STOP = s causes STEPWISE to stop when it has
found the "best" s-variable model,
where s is the STOP value. This
applies only to the MAXR or MINR
methods.

(3) Sample

```
//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SYSIN DD *
DATA FITNESS;
INPUT AGE WEIGHT OXY RUNTIME RSTPULSE
RUNPULSE MAXPULSE;
CARDS;
44 89.47 44.609 11.37 62 178 182 40 75.07 45.313 10.07 62 185 185
44 85.84 54.297 8.65 45 156 168 42 68.15 59.571 8.17 40 166 172
38 89.02 49.874 9.22 55 178 180 47 77.45 44.811 11.63 58 176 176
40 75.98 45.681 11.95 70 176 180 43 81.19 49.091 10.85 64 162 170
44 81.42 39.442 18.08 63 174 176 38 81.87 60.055 6.63 48 170 186
44 73.03 50.541 10.13 45 168 168 45 87.66 37.388 14.03 56 186 192
45 66.45 44.754 11.12 51 176 176 47 79.15 47.273 10.60 47 162 164
54 83.12 51.855 10.33 50 166 170 49 81.42 49.156 8.95 44 180 185
51 69.63 40.836 10.95 57 168 172 51 77.91 46.672 10.00 48 162 168
48 91.63 46.774 10.25 48 162 164 49 73.37 50.388 10.08 67 168 168
57 73.37 39.407 12.63 58 174 176 54 79.38 46.080 11.17 62 156 165
52 76.32 45.441 9.63 48 164 166 50 70.87 54.625 8.92 48 146 155
51 67.25 45.118 11.08 48 172 172 54 91.63 39.203 12.88 44 168 172
51 73.71 45.790 10.47 59 186 188 57 59.08 50.545 9.93 49 148 155
49 76.32 48.673 9.40 56 186 188 48 61.24 47.920 11.50 52 170 176
52 82.78 47.467 10.50 53 170 172
;
PROC STEPWISE;
MODEL OXY=RUNTIME AGE WEIGHT RUNPULSE
MAXPULSE RSTPULSE/FORWARD BACKWARD MAXR;
/*
//
```

SAS

FORWARD SELECTION PROCEDURE FOR DEPENDENT VARIABLE GRV

STEP 1 VARIABLE HUNTING ENTERED R SQUARE = 0.51278375 CIP1 = 8.66735310

	DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	3	135.46619793	135.46619793	14.62	0.0019
ERROR	14	129.74622900	9.26758779		
TOTAL	15	265.21242694			
B VALUE STD ERROR TYPE II SS F PROB>F					
INTERCEPT	61.64683045				
HUNTING	-1.39644276	0.36877384	135.46619793	14.62	0.0019

STEP 2 VARIABLE WEIGHT ENTERED R SQUARE = 0.60399406 CIP1 = 5.11058824

	DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	2	160.18726178	80.09363089	9.91	0.0024
ERROR	13	105.02516516	8.07658206		
TOTAL	15	265.21242694			
B VALUE STD ERROR TYPE II SS F PROB>F					
INTERCEPT	46.20605094				
HUNTING	-1.13948138	0.34831710	123.07815246	15.23	0.0018
WEIGHT	0.16201968	0.09782092	24.72106384	3.06	0.1098

STEP 3 VARIABLE AGE ENTERED R SQUARE = 0.64483145 CIP1 = 5.47608950

	DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	3	171.0473904	57.18246348	7.33	0.0047
ERROR	12	97.06448790	7.40539066		
TOTAL	15	265.21242694			
B VALUE STD ERROR TYPE II SS F PROB>F					
INTERCEPT	59.95801942				
HUNTING	-1.37103355	0.33932785	128.17854938	16.42	0.0016
AGE	-0.17943680	0.14641860	11.36018726	1.46	0.2500
WEIGHT	0.12445818	0.09421845	13.05990043	1.67	0.2202

STEP 4 VARIABLE RSTPULSE ENTERED R SQUARE = 0.7123553 CIP1 = 4.76887076

	DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	4	197.22507929	49.30626982	6.98	0.0048
ERROR	11	76.98734764	6.99885140		
TOTAL	15	265.21242694			
B VALUE STD ERROR TYPE II SS F PROB>F					
INTERCEPT	64.54755177				
HUNTING	-1.02647845	0.37895301	30.12552050	7.35	0.0202
AGE	-0.23048021	0.14810594	18.79815164	2.75	0.1259
WEIGHT	0.12516517	0.09992202	12.78907088	1.80	0.1981
RSTPULSE	-0.19485853	0.11712252	18.67734026	2.74	0.1261

STEP 5 VARIABLE MAXPULSE ENTERED R SQUARE = 0.73107392 CIP1 = 6.26159313

	DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	5	175.88988811	35.17797766	5.44	0.0115
ERROR	10	71.32253882	7.13225388		
TOTAL	15	245.21242694			
B VALUE STD ERROR TYPE II SS F PROB>F					
INTERCEPT	52.27960122				
HUNTING	-0.89271791	0.42982634	30.76615923	4.31	0.0445
AGE	-0.22558044	0.14872712	17.59575919	2.48	0.1477
WEIGHT	0.14002908	0.09493859	15.51592637	2.18	0.1710
MAXPULSE	0.10425073	0.14544037	3.46640392	0.51	0.4899
RSTPULSE	-0.28827682	0.17428374	19.02444138	2.64	0.1253

NO OTHER VARIABLES MET THE 0.5000 SIGNIFICANCE LEVEL FOR ENTRY INTO THE MODEL.

SAS
FORWARD SELECTION PROCEDURE FOR DEPENDENT VARIABLE DRY

STEP 4		VARIABLE PUMPULSE ENTERED	R SQUARE = 0.76413655	CPI = 7.00000000		
		DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	6	202.65859068	33.77641828	4.86	0.0175	
ERROR	9	62.55391725	6.95035225			
TOTAL	15	265.21242094				
		S VALUE	STD ERROR	TYPE II SS	F	PROB>F
INTERCEPT	49.36009962	0.44342815	19.78300062	2.85	0.1259	
RUNTIME	-0.74804863	0.14203193	18.47769907	2.37	0.1980	
AGE	-0.21884196	0.09925160	7.93197871	1.08	0.3250	
WEIGHT	0.10239611	0.27958750	8.76862137	1.26	0.2904	
PUMPULSE	-0.31643653	0.34350635	13.18251832	1.79	0.2182	
MAXPULSE	-0.09642012	0.17950198	23.73503169	3.41	0.0977	
RSTPULSE	-0.33238164					

STEP 0		ALL VARIABLES ENTERED	R SQUARE = 0.76413655	CPI = 7.00000000		
		DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	6	202.65859068	33.77641828	4.86	0.0175	
ERROR	9	62.55391725	6.95035225			
TOTAL	15	265.21242094				
		S VALUE	STD ERROR	TYPE II SS	F	PROB>F
INTERCEPT	49.36009962	0.44342815	19.78300062	2.85	0.1259	
RUNTIME	-0.74804863	0.14203193	18.47769907	2.37	0.1980	
AGE	-0.21884196	0.09925160	7.93197871	1.08	0.3250	
WEIGHT	0.10239611	0.27958750	8.76862137	1.26	0.2904	
PUMPULSE	-0.31643653	0.34350635	13.18251832	1.79	0.2182	
MAXPULSE	-0.09642012	0.17950198	23.73503169	3.41	0.0977	
RSTPULSE	-0.33238164					

STEP 1		VARIABLE HEIGHT REMOVED	R SQUARE = 0.75572922	CPI = 6.08395783		
		DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	5	195.12853507	39.02490616	5.57	0.0104	
ERROR	10	70.08388586	7.00838859			
TOTAL	15	265.21242094				
		S VALUE	STD ERROR	TYPE II SS	F	PROB>F
INTERCEPT	62.01457864	0.44461896	21.78891350	3.11	0.1086	
RUNTIME	-0.78284870	0.13649729	25.72059967	2.87	0.1084	
AGE	-0.26148283	0.08131365	18.78054903	2.59	0.1532	
PUMPULSE	-0.40084066	0.33778480	17.08377933	2.43	0.1501	
MAXPULSE	0.52689645	0.17828524	21.93326279	3.07	0.1103	
RSTPULSE	-0.31410622					

STEP 2		VARIABLE PUMPULSE REMOVED	R SQUARE = 0.67247015	CPI = 6.49396072		
		DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	4	178.37396196	44.59349049	5.65	0.0101	
ERROR	11	86.83845899	7.89440591			
TOTAL	15	265.21242094				
		S VALUE	STD ERROR	TYPE II SS	F	PROB>F
INTERCEPT	78.67800249	0.44603645	41.06128972	5.20	0.0435	
RUNTIME	-1.01284859	0.14316093	33.23717104	4.21	0.0448	
AGE	-0.20374820	0.14024642	9.93750333	0.12	0.7368	
MAXPULSE	0.05159850	0.18774957	13.70845793	1.74	0.2144	
RSTPULSE	-0.24213663					

STEP 3		VARIABLE MAXPULSE REMOVED	R SQUARE = 0.66933554	CPI = 4.62890961		
		DF	SUM OF SQUARES	MEAN SQUARE	F	PROB>F
REGRESSION	3	177.63600862	59.14533614	8.09	0.0033	
ERROR	12	87.57641232	7.31470154			
TOTAL	15	265.21242094				
		S VALUE	STD ERROR	TYPE II SS	F	PROB>F
INTERCEPT	82.73561050	0.39040747	56.48867563	7.58	0.0175	
RUNTIME	-1.07495012	0.13779178	35.09228680	4.56	0.0546	
AGE	-0.29328149	0.11199384	16.96823963	2.59	0.1535	
RSTPULSE	-0.19629927					

ALL VARIABLES IN THE MODEL ARE SIGNIFICANT AT THE 0.1000 LEVEL.

B. 분산분석 (Analysis of Variance)

1. 분산분석

- (1) 분산분석은 관측치 사이의 분산을 각 요인에 따른 분산으로 나누는 방법이다. 분산분석의 방법은 실험 (Experiment)이 어떻게 계획되었는가에 따라 다르며 여기에서는 Completely Randomized Design, Randomized Blocks Design에 대하여 설명하기로 한다.

Completely Randomized Design (CRD)

- (2) 모형 $y_{ij} = \mu + \tau_i + \epsilon_{ij}$
- $\begin{array}{ccccccc} \uparrow & & \uparrow & & \uparrow & & \uparrow \\ \text{반응} & & \text{평균} & & \text{처리 } i \text{의 효과} & & \text{오차} \\ \text{(response)} & & & & \left(\sum_{i=1}^t \tau_i = 0 \right) & & \end{array}$

- (3) CRD는 t 종류의 처리 (treatment)의 효과를 알아보기 위한 실험 계획이며 각 처리에 n개의 subject를 임의로 배치한다.

2. ANOVA procedure.

(1) PROC ANOVA;

```
CLASS treatment variable;  
MODEL response = treatment variable;
```

(2) SAS EXAMPLE

PROC ANOVA;

CLASS BRAND;

MODEL WEAR = BRAND;

MEANS BRAND/DUNCAN;

```
//SPCSAS30 JOB CLASS=V  
//SAS EXEC SAS  
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR  
//SYSIN DD *  
DATA;  
INPUT BRAND* WEAR;  
CARDS;  
ACME 2.3  
ACME 2.4  
ACME 2.1  
ACME 2.5  
CHAMP 2.2  
CHAMP 2.3  
CHAMP 2.4  
CHAMP 2.6  
AJAXP 2.2  
AJAXP 2.0  
AJAXP 1.9  
AJAXP 2.1  
TUFFY 2.4  
TUFFY 2.7  
TUFFY 2.6  
TUFFY 2.8  
XTRA 2.4  
XTRA 2.5  
XTRA 2.3  
XTRA 2.4  
;  
PROC PRINT;  
PROC ANOVA;  
CLASS BRAND;  
MODEL WEAR=BRAND;  
MEANS BRAND/DUNCAN;  
;  
/*  
//
```

17:41 THURSDAY, NOVEMBER 15, 1984 1

```

      SAS
OBS   BRAND   WEAR
  1   ACME    2.3
  2   ACME    2.4
  3   ACME    2.1
  4   ACME    2.5
  5   CHAMP   2.2
  6   CHAMP   2.3
  7   CHAMP   2.4
  8   CHAMP   2.6
  9   AJAXP   2.2
 10  AJAXP   2.0
 11  AJAXP   1.9
 12  AJAXP   2.1
 13  TUFFY    2.4
 14  TUFFY    2.7
 15  TUFFY    2.6
 16  TUFFY    2.8
 17  XTRA     2.4
 18  XTRA     2.5
 19  XTRA     2.3
 20  XTRA     2.4

```

17:41 THURSDAY, NOVEMBER 15, 1984 2

```

      SAS
ANALYSIS OF VARIANCE PROCEDURE
      CLASS LEVEL INFORMATION
CLASS   LEVELS   VALUES
BRAND   5        ACME AJAXP CHAMP TUFFY XTRA

```

NUMBER OF OBSERVATIONS IN DATA SET = 20

SAS

BT141 THURSDAY, NOVEMBER 15, 1984 3

ANALYSIS OF VARIANCE PROCEDURE

DEPENDENT VARIABLE: MEAR

SOURCE	DF	SUM OF SQUARES	MEAN SQUARE	F VALUE	PR > F	R-SQUARE	C-Y.
MODEL	4	0.4770000	0.11925000	7.64	0.0015	0.070629	0.3221
ERROR	15	0.33250000	0.02216667		ADJUSTED MSE		MEAN MEAN
CORRECTED TOTAL	19	1.00950000			0.14680676		2.85500000

SOURCE	DF	ANOVA SS	F VALUE	PR > F
BRAND	4	0.47700000	7.64	0.0015

SAS

ANALYSIS OF VARIANCE PROCEDURE

DUNCAN'S MULTIPLE RANGE TEST FOR VARIABLE: MEAR
 NOTE: THIS TEST CONTROLS THE TYPE I COMPARISONWISE ERROR RATE,
 NOT THE EXPERIMENTWISE ERROR RATE.

ALPHA=0.05 DF=15 MSE=.0221667

MEANS WITH THE SAME LETTER ARE NOT SIGNIFICANTLY DIFFERENT.

DUNCAN	GROUPING	MEAN	N	BRAND
A		2.6250	4	TUFFY
B		2.4000	4	XTRA
B		2.3750	4	CHAMP
B		2.3250	4	ACHE
C		2.0500	4	AJAXP

SAS

14:41 WEDNESDAY, OCTOBER 31, 1984 1

OBS	BLOCK	BLEND	PCTLOSS
1	1	B	18.2
2	1	A	19.6
3	1	C	17.0
4	1	E	18.3
5	1	D	15.1
6	2	A	16.5
7	2	E	16.3
8	2	B	19.2
9	2	C	18.1
10	2	D	16.0
11	3	B	17.1
12	3	D	17.8
13	3	C	17.3
14	3	E	15.8
15	3	A	17.5

SAS

14:41 WEDNESDAY, OCTOBER 31, 1984 2

ANALYSIS OF VARIANCE PROCEDURE

CLASS LEVEL INFORMATION

CLASS	LEVELS	VALUES
BLOCK	3	1 2 3
BLEND	5	A B C D E

NUMBER OF OBSERVATIONS IN DATA SET = 15

SAS

14:41 WEDNESDAY, OCTOBER 31, 1984 3

ANALYSIS OF VARIANCE PROCEDURE

DEPENDENT VARIABLE:	DF	SUM OF SQUARES	MEAN SQUARE	F VALUE	PR > F	R-SQUARE	C.V.
SOURCE							
MODEL	4	10.66100000	2.66525000	1.10	0.4370	0.452260	7.1650
ERROR	4	12.39400000	3.09850000				
CORRECTED TOTAL	14	23.05500000			1.2694502		17.3200000
SOURCE	DF	MEAN SS	F VALUE	PR > F			
BLOCK	2	5.34000000	0.80	0.3778			
BLEND	4	10.49400000	1.61	0.2615			

(4) Example

```
PROC ANOVA;
```

```
CLASS BLOCK BLEND;
```

```
MODEL PCTLOSS = BLOCK BLEND;
```

```
MEANS BLEND/WALLER;
```

Note: PCTLOSS = percent loss of insects

BLOCK = locations (b=3)

BLEND = blends of household insecticide (t-s)

C. t - 검정 (t - Test)

1. Two related samples

- (1) 두 측정치가 쌍 (pair) 을 이루고 있는 경우의 t = test 를 말한다. 예를 들면 실험의 subject 로 15마리의 동물이 채택되고 각 동물이 흥분안정제를 주입받는다고 할 때 투약전의 맥박수를 x , 투약후의 맥박수를 y 라 하면 (x , y) 는 쌍을 이룬다. x 와 y 의 값에 차이가 있는지의 여부를 알기 위한 t statistic 은

$$t = \frac{\bar{d}}{s\bar{d}/n}$$

이고 여기에서 $\bar{d} = \bar{y} - \bar{x}$, $s\bar{d}^2 = \frac{\sum_{i=1}^n (d_i - \bar{d})^2}{(n-1)}$,

$$d_i = y_i - x_i \quad (i=1, \dots, n) \text{ 이다.}$$

- (2) SAS : Example

```
//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SYSIN DD *
DATA SAMPLE;
  INPUT ID PRETEST POSTTEST;
  DIFF=POSTTEST-PRETEST;
CARDS;
1 80 82
2 73 71
3 70 95
4 60 69
5 88 100
6 84 71
7 65 75
8 37 60
9 94 95
10 98 99
11 52 65
12 78 83
13 40 60
14 79 86
15 59 62
;
PROC PRINT;
PROC MEANS MEAN STDERR T PRT;
  VAR DIFF;
TITLE PAIRED-COMPARISONS T TEST;
/*
//
```


17:43 THURSDAY, NOVEMBER 15, 1984 1

SAS

OBS	ID	PRETEST	POSTTEST	DIFF
1	1	80	82	2
2	2	73	71	-2
3	3	70	95	25
4	4	60	69	9
5	5	88	100	12
6	6	84	71	-13
7	7	65	75	10
8	8	37	60	23
9	9	91	95	4
10	10	98	99	1
11	11	52	65	13
12	12	78	83	5
13	13	40	60	20
14	14	79	86	7
15	15	59	62	3

17:43 THURSDAY, NOVEMBER 15, 1984 2

PAIRED-COMPARISONS T TEST

VARIABLE	MEAN	STD ERROR OF MEAN	T	PR> T
DIFF	7.9333333	2.56434651	3.09	0.0079

2. Two Independent Samples

- (1) 두 표본이 독립적으로 얻어진 경우 평균치간의 차이가 있는지의 여부를 알아보는 t - test 이다. 첫번째 표본을 (x_1, \dots, x_{n_1}) 이라고 하고 두번째의 표본을 (y_1, \dots, y_{n_2}) 이라고 하면

$$t = (\bar{y} - \bar{x}) \sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

이고 여기에서 $s^2 = \left[\sum_{i=1}^{n_1} (x_i - \bar{x})^2 + \sum_{i=1}^{n_2} (y_i - \bar{y})^2 \right] / (n_1 + n_2 - 2)$ 이다.

- (2) 예를 들면 두 종류의 탄약으로 만들어진 총알의 속도를 비교해 보는 경우 첫번째 표본은 한 종류의 탄약으로 만들어진 총알의 속도 측정치로 이루어져 있고 두번째 표본은 다른 종류의 탄약으로 만들어진 총알의 속도 측정치로 이루어져 있다고 하자.
- (3) SAS : Example

```
//SPCSAS30 JOB CLASS=V
//SAS EXEC SAS
//SAMPLE DD DSN=SAS.SAMPLE,DISP=SHR
//SYSIN DD *
DATA;
INPUT POWDER VELOC;
CARDS;
1 27.3
1 28.1
1 27.4
1 27.7
1 28.0
1 27.4
1 27.1
1 28.1
2 28.3
2 27.9
2 28.2
2 28.4
2 27.9
2 27.7
2 28.5
2 27.9
2 28.4
2 27.8
;
PROC PRINT;
PROC TTEST;
CLASS POWDER;
TITLE T-TEST SAMPLE;
```

14:40 WEDNESDAY, OCTOBER 31, 1984 1

SAS		
OBS	PGWDR	VELDC
1	1	27.3
2	1	28.1
3	1	27.4
4	1	27.7
5	1	28.0
6	1	27.4
7	1	27.1
8	1	28.1
9	2	28.3
10	2	27.9
11	2	28.2
12	2	28.4
13	2	27.9
14	2	27.7
15	2	28.5
16	2	27.9
17	2	28.4
18	2	27.8

T-TEST SAMPLE
TTEST PROCEDURE

14:40 WEDNESDAY, OCTOBER 31, 1984 2

VARIABLE: VELDC

PGWDR	N	MEAN	STD DEV	STD ERROR	MINIMUM	MAXIMUM	VARIANCES	F	DF	PROB > F
1	9	27.63778	0.38256482	0.12677765	27.10000000	28.10000000	UNEQUAL	-2.7785	12-8	0.0165
2	10	28.1200000	0.29059326	0.09189366	27.70000000	28.50000000	EQUAL	-1.8762	8-8	0.0810

FOR MEAN VARIANCES ARE EQUAL, F= 1.82 WITH T AND 9 DF PROB > F= 0.1946

CHAPTER 5. 다 변 량 분 석

CHAPTER 5. 다변량분석

A. DISCRIMINANT Analysis

1. Introduction

- (1) Data: Multivariate (correlated) numerical data with pre-defined group identifier (class variable).
- (2) Object:
 1. Develop a decision rule (discriminant model) which can separate maximally based on the information of the sample.
 2. Assign some new observations with unknown origin into one of the known groups.

2. Related Procedures

- (1) DISCRIM: For the approximate multivariate normal population within group (with equal or not-equal variance-covariance matrix).
- (2) NEIGHBOR: For radically non-normal population using non-parametric nearest neighbor method.

- (3) CANDISC: Dimension-reduction technique related principal component and canonical correlation by taking linear composites of response variables with class variable, which give maximum between group deviation.
- (4) STEPDISC: As a variable selection technique using forward selection, backward elimination or stepwise selection method to find a subset of variables that best discriminates group differences.

3. Alternative Procedure

- (1) FUNCAT: Fitting categorical linear model with classification variable as the dependent variable.
- (2) ANOVA: Series of Univariate Analysis of variance technique.

4. Background of DISCRIM Procedure

- (1) Distribution in each group approx multivariate normal

$$N(\underline{\mu}^t, \Sigma^t) \text{ or } N(\underline{\mu}^t, \Sigma)$$

for $t=1,2,\dots,k$ (number of groups).

μ_t	Z_t	Σ	:	Population parameter
\updownarrow	\updownarrow	\updownarrow		
M_t	S_t	S	:	Sample estimator

(2) Homogeneity test on Σ_t ($t=1,2,\dots,k$) : Use the Bartlett's likelihood ratio test Statistic (Approximate Chi-square test) and the results of test determine whether the criterion is based on S_t or S .

(3) Decision Rule: Generalized Squared distance from \underline{x} to group t such that

$$D_t^2(\underline{x}) = g_1(\underline{x}, t) + g_2(t), \text{ where}$$

$$g_1(\underline{x}, t) = (\underline{x} - \underline{m}_t)' S_t^{-1} (\underline{x} - \underline{m}_t) + \log_e |S_t|$$

or

$$(\underline{x} - \underline{m}_t)' S (\underline{x} - \underline{m}_t)$$

"Sample Mahalanobis distance"

and $g_2(t) = g - 2 \log_e q_t$, where

q_t = Prior Probability of group t

If q_t = constant for all t then $q_t = 0$.

Assign a new \underline{x} to group t if

$$D_t^2(\underline{x}) = \min(D_1^2, D_2^2, \dots, D_k^2)$$

Equivalently SAS uses posterior probability of \underline{x} belonging to group u, where

Posterior Probability

$$P_u(\underline{x}) = e^{-\frac{1}{2}D_u(\underline{x})} / \sum_{i=1}^k e^{-\frac{1}{2}D_i(\underline{x})}$$

and an observation X is assigned to group u if

$$P_u(\underline{x}) = \max\{P_1(x), \dots, P_k(x)\}$$

5. Outline of Use

```
PROC DISCRIM DATA = XX SIMPLE POOL = TEST W CORR OUT =  
OUTXX; CLASS GROUP;  
PROC DISCRIM DATA = OUTXX TESTDATA = YY;  
TESTCLASS GROUP;
```

B. FACTOR Analysis

1. Introduction

- (1) Data: Multivariate Numerical data, Correlation matrix, Covariance matrix, Factor Pattern matrix, or Scoring coefficient matrix.
- (2) Object: Investigate structural relationship among response variables.

$\left\{ \begin{array}{l} \text{Several,} \\ \text{Difficult to interpret,} \\ \text{Correlated} \end{array} \right\}$	variables
--	-----------

Using common factor	$\left\{ \begin{array}{l} \text{Few,} \\ \text{Conceptually meaningful,} \\ \text{Relatively independent} \end{array} \right\}$	Hidden factors generating the dependence in responses
---------------------------	---	--

analysis with rotation.

2. Related Procedures (in the sense of structural analysis)

- (1) CANCELL: Relationship between two sets of variables by finding a small number of linear composites for each set (canonical variables).

- (2) PRINCOMP: Obtain a small number of linear composites of original variables (Principal Components) by orthogonal transformation that has as much the information in the original variables as possible (variance-oriented).

- (3) FACTOR:

3. Background of FACTOR Procedure

Model: $x = \lambda_{u1}f_1 + \lambda_{u2}f_2 + \dots + \lambda_{uk}f_k + e^i, i=1,2,\dots, p$ (# of responses) and $P \gg k$

, where

$f_1, f_2, \dots, f_k : k$ - common-factor variates

(factor loadings) : parameters representing the importance of the j 'th factor on i 'th response (need not to be orthogonal)

$e^{(i)}$: i 'th specific factor variate.

$x^{(i)}$: linear function of a small number of unobservable (hidden) latent, common factor)) variables and a single (latent) specific factor. (i.e., response variable).

$$\underline{x} = A\underline{f} + \underline{e} \quad , \text{ where}$$

$f = N(Q, I), \underline{e} = N(0, D\varphi_i)$ φ_i is i 'th specific variance (specificity) and \underline{y} and \underline{e} are stochastically independent.

Then the fundamental representation of Z under factor model is

$$\Sigma = AA' + D\varphi_i$$

Thus

$$\text{Var}(x^{(i)}) = \sigma_i^2 - \varphi_i = \sum_{j=1}^k \lambda_{ij}^2 \quad (\text{communality of the } i\text{'th response})$$

$$\text{Cov}(x^{(i)}, x^{(j)}) = \sigma_{ij} = \sum_{l=1}^k \lambda_{il} \lambda_{jl}$$

and $\text{cov}(\underline{x}, \underline{f}') = A$ (or correlation of \underline{x} and \underline{f} if standardized)

4. Factor Rotation

Let $AT = A^*$ such that $T T' = I$ (orthogonal) then $\Sigma = A^* A^* = D\psi = M + D\psi$. Thus there exists infinitely many such T (i.e., infinitely many factor loading matrix generating the same Σ).

Thurstone (1945) proposed the concept of "simple structure" as a means of selecting A most meaningful and interpretable (simplicity and parsimony).

Orthogonal rotation	}	Varimax, quartimax, equamax, procrustean
Oblique rotation		

e.g., Varimax (Kaiser) : Minimize the sum of variances of the squared loadings within each column of factor matrix.

5. Factor Extraction and Goodness of Fit Test on the Factor Model

(1) Maximum likelihood factor analysis:

Factor extraction by maximum likelihood technique, and goodness of fit test for the number of factors extracted by the generalized likelihood ratio principle. (Lawley and HOME).

(2) Principal component analysis.

(3) Principal factor analysis.

(4) Iterative Principal Factor analysis.

(5) Alpha factor analysis.

(6) Image analysis.

6. Outline of Use

(1) Principal Component Analysis:

```
PROC FACTOR SCORE OUTSTAT = XX
  SCREE REORDER;
PROC SCORE DATA = A SCORE =
  OUT = SCDATA;
PROC PLOT; PLOT FACTOR 2 * FACTOR 1
```

- (2) Principal factor analysis (Simplest and Computationally efficient) :

PROC FACTOR ... ; PRIORS $\begin{pmatrix} \text{SMC} \\ \text{NAX} \end{pmatrix}$; (for singular correlation matrix)

PROC FACTOR N = #

ROTATE = VARIMAX ROUND SCORE ... ;

PROC SCORE ... ;

PROC PLOT ... ;

PROC REG ... ;

- (3) Maximum likelihood factor analysis: (best in statistical point of view)

- 1) Rigid test for the number of factors.
- 2) Desirable asymptotic properties on the estimators.
- 3) Not require normality but expensive.

PROC FACTOR METHOD = ML;

C. CLUSTER Analysis

1. Introduction

- (1) Data: 1. Multivariate numerical data without pre-defined group information.

2. Similarity matrix (squared distance matrix, correlation type matrix).
- (2) Object: Assign (cluster) observations (or variable) into groups in such a way that (high) similarity within group and (high) dissimilarity between groups in some sense.
 - (3) Types of clusters
 1. Disjoint cluster
 2. Hierarchical cluster
 3. Overlapping cluster
 4. Fuzzy cluster (combined)

2. Related Procedures

- (1) CLUSTER: Find hierarchical clusters of observations using centroid, Ward-Hook, or average linkage method on squared Euclidean distance.
- (2) FASTCLUS: Disjoint cluster based on K-means method for large data.
- (3) VARCLUS: Hierarchical and disjoint clustering of variables.

- (4) TREE: Drew a tree diagram (dendrogram) using output from CLUSTER or VARCLUS.

3. CLUSTER Procedure

Use the three standard agglomerative hierarchical clustering algorithms based on different distance measures as a measure of similarity.

- (1) Centroid : Euclidean distance between cluster centroids
(Robust to outliers).
- (2) Ward-Hook Method : Error sum of squares within clusters,
within clusters,

$$\begin{aligned}
 x_{ij} & \quad k \quad (i = 1, 2, \dots, p \text{ (\# of responses)}) \\
 & \quad j = 1, 2, \dots, n_k \text{ (size of } k\text{th cluster)} \\
 & \quad k = 1, 2, \dots, c \text{ (\# of clusters)}.
 \end{aligned}$$

$$(\text{SSE})_k = \sum_{i=1}^p \sum_{j=1}^{n_k} (x_{ijk} - \bar{x}_{ik})^2$$

$$\text{SSE} = \sum_{k=1}^c (\text{SSE})_k$$