

# 関西SASユーザ-会

## -SASを使った便利なTIPS&テクニック-

SAS Institute Japan Ltd.  
東 一成 (Kazunari Azuma)  
Kazunari.Azuma@sas.com

*The Power to Know.*

## 本日のアジェンダ

- ◆ あるユーザーからの質問からサンプルを二つご紹介します。  
(BaseSASだけではなくSAS/STATも必要ですが...)
  - SASでランダムサンプリングを行うにはどうすればよいですが？
    - ◆ 色々な手法を選んだり、件数や繰り返してサンプリングを行う簡単な方法はありませんか？
  - 主成分分析により、Prin1、Prin2、Prin3、...と変数が作成されます。Eigenvalueが1以上の変数のみをKeepするには？

## SASによるサンプリングについて

## サンプリングのための方法

- ◆ V6の場合・・・
  - RANUNI関数を利用して、乱数を発生させてデータステップなどを利用して抽出する。
  - SQLプロシジャなどを用いて抽出する。
  - 層別抽出の場合はFREQプロシジャなどを組み合わせて実現をする。
  - BaseSASがあればOK。
- ◆ V8の場合・・・
  - SURVEYSELECTプロシジャを利用すれば数行で済む。
  - ただしSAS/STATが必要。

## 色々なサンプリング手法

- ◆ 例えば以下のようなサンプリング手法を考えます。
  - 単純無作為抽出
  - 系統抽出法(システムサンプリング)
  - 層別抽出
  - 復元抽出
  - 非復元抽出
  
- 上記の抽出方法をいかにして行うかをご紹介します。
  
- ◆ V6による方法は単純無作為抽出と層別抽出の二つだけにしたいと思います。

5

## それぞれのサンプリング手法について1

- ◆ 単純無作為抽出法
  - 単純無作為抽出とは、母集団のどの個体も同じ確率で取り出される抽出法。乱数表やコンピューターに乱数を発生をさせて、サンプルを抽出する場合が多い。
- ◆ 系統抽出法(システムサンプリング)
  - 系統抽出法は母集団構成員のリストで標本を等間隔に並んでいる個体を取り出していく方法。500人の母集団から50人のサンプルが必要だとすると、最初の1名は乱数などにより決定される。最初が5番だとするとそれ以降は、15、25、35...といった形で抽出されて行く方法。
- ◆ 層別抽出法
  - 母集団をいくつかの層に分け、それぞれの層から単純無作為抽出法などによって、サンプルを作成する方法。各層から抽出する場合、母集団の各層の人数に比例させる方法が多い。これは各母集団の内部で、いくつかの均質な層に分かれていると考え、各層別に標本を抽出してサンプルを作成した方が誤差が小さく、均質なサンプルを得ることができるという考え方。

6

## それぞれのサンプリング手法について2

- ◆ 復元抽出
  - 母集団から標本を抽出する場合に、どの標本も同じ確率で取り出される状態である必要がある。このためには取り出した標本をもう一度母集団に戻して、次の標本を取り出す必要がある。ただ場合によっては重複されて抽出される標本が出てくる。
- ◆ 非復元抽出
  - 上記の復元抽出とは逆で、一度抽出された標本は母集団に戻さず、次の標本を抽出する。
- ◆ 上記の二つの手法であるが、サンプルの標本数に比較して、母集団が十分に大きければ、重複が起こる可能性も低いので、差はほとんど無いと考えられる。

## 今回利用するデータについて

- ◆ 今回のサンプリングの例で利用するデータは以下のような形で、20,000件あります。

OBS	CUST_ID	AGE	GENDER	FLG	Sales	MARRIAGE	HOUSING
1	1	44	女性	0	1500	既婚	学生寮
2	2	44	女性	1	0	未婚	学生寮
3	3	37	女性	0	0	既婚	社宅
4	4	56	女性	0	0	既婚	マンション(賃)
5	5	32	女性	0	3400	既婚	社宅
6	6	62	女性	0	0	既婚	マンション(賃)
7	7	60	女性	0	0	未婚	マンション(賃)
8	8	40	男性	1	0	未婚	社宅
9	9	51	女性	1	0	未婚	社宅
10	10	27	女性	0	0	未婚	学生寮
11	11	44	女性	1	0	未婚	学生寮
12	12	28	不明	0	12000	既婚	マンション(持)
13	13	31	女性	0	0	既婚	社宅
14	14	34	女性	1	0	既婚	マンション(持)

## 無作為抽出の方法-SASV6の場合 1

```
/*件数を指定する場合-SQLプロシジャ*/
proc sql outobs=2000;
  create table sample1 as
  select ranuni(0) as random, * from ksug.customer order
  by random;
quit;
```

SQLプロシジャを利用して行う方法。この場合はRANUNI関数を利用して乱数を発生させ、乱数の値で上位2000件を取ってきている。比較的単純に行う事ができるので便利。

9

## 無作為抽出の方法-SASV6の場合 2

```
/*件数を指定する場合-データステップ*/
data sample2 ;
  set ksug.customer nobs = total ;
  if _sample_count_ < 2000 then do ;
    if ranuni(0) *(total + 1 - _N_) <= ( 2000 - _sample_count_) then do ;
      _sample_count_ + 1 ;
      output ;
    end ;
  end ;
  drop _sample_count_ ;
run ;
```

データステップを使って行う方法は上記の通り。乱数を発生させる部分はSQLプロシジャの場合と同じであるが、上記の例では2,000件と指定してあるので、備ったとり方をしないように工夫が行われている。流れとしては以下のとおりになる。

例えば・・・20,000件のうち2,000件を抽出したいので、全体の約10%を抽出。  
 1件目 乱数0.5 \* 20000+1 - 1 <= (2000 - 0) 抽出されない  
 2件目 乱数0.1 \* 20000+1 - 2 <= (2000 - 0) 抽出される  
 3件目 乱数0.9 \* 20000+1 - 3 <= (2000 - 1) 抽出されない  
 上記の処理を2000件になるまで行う。

10

## 無作為抽出の方法-SASV6の場合 3

```
/*おおよその割合で抽出する場合*/
data sample3;
  set ksug.customer;
  if ranuni(0)<=0.1;
run;
```

上記の場合も乱数を発生させる方法と同じですが、件数の指定をしていないので記述が簡単。ただ、20,000件のうち2,000件を抽出したいと思っても、数件~数十件程度のプレが見られる。

## 層別抽出の方法-SAS V6での場合

```
/*性別の[男性][女性]の各層から10%ずつ抽出*/
data sample4(drop=n1--_cnt2);
set ksug.customer;
if gender='男性' then do;
  n1+1;
  if _cnt1_ < 820 then do ;
    if ranuni(0) *(8205 - n1) <= ( 820 - _cnt1_) then do ;
      _cnt1_ + 1 ;output ;
    end ;
  end ;
end ;
if gender='女性' then do;
  n2+1;
  if _cnt2_ < 807 then do ;
    if ranuni(0) *(8069 - n2) <= ( 807 - _cnt2_) then do ;
      _cnt2_ + 1 ;output ;
    end ;
  end ;
end ;
run;
```

データステップなどを用いて層別抽出を行うには、左の図のようになる。

性別の各層から10%ずつ抽出する処理であるが、まずFREQプロシジャなどで各層の件数を把握し、その内10%の件数を抽出できるようにロジックを組み合わせていく必要がある。

## V8以降(SAS/STAT必要)でのサンプリング

- ◆ SAS V8からは新たにSURVEYSELECTプロシジャが追加され、様々なサンプリング手法を行う事ができるようになっている。
- ◆ 構文は以下のとおり。  
PROC SURVEYSELECT options ;  
STRATA variables ;  
CONTROL variables ;  
SIZE variable ;  
ID variables ;
- ◆ 手法は11通り用意されているが、今回は以下の手法を紹介します。
  - 単純無作為抽出
  - システムサンプリング(系統抽出法)
  - 層別抽出
  - 復元抽出
  - 非復元抽出

## 単純無作為抽出 1

```
/*件数を指定して行う非復元抽出*/
PROC SURVEYSELECT DATA=ksug.customer
METHOD=SRS
N=2000
OUT=sample5;
RUN;QUIT;
```

最も単純なサンプリングを行っている。手法は「METHOD=SRS」で指定している。SRSはSimple Random Samplingの事。

結果は指定した通りに2000件のデータセットがサンプルとして出力。

観測番号	年齢	性別	収入	職業	婚姻状況	投票意向
1	35	30 女性	0	18728 既婚	投票	
2	29	40 不明	0	8 未婚	アバート	
3	30	40 女性	0	8 未婚	投票	
4	31	35 女性	0	4608 未婚	一戸建て	
5	38	27 女性	0	41696 既婚	マンション(独)	
6	30	40 女性	0	8 既婚	マンション(独)	
7	44	39 女性	0	8 既婚	住宅	
8	66	34 女性	0	8 未婚	アバート	
9	32	42 女性	0	8 未婚	マンション(独)	
10	60	39 女性	0	22008 既婚	マンション(独)	
11	50	36 女性	0	8 未婚	マンション(独)	
12	111	39 女性	0	8426 未婚	アバート	
13	136	38 女性	8	90908 既婚	住宅	
14	132	37 不明	0	75214 既婚	学生寮	
15	127	30 女性	0	8 既婚	マンション(独)	

## 単純無作為抽出 2

**/\*件数を指定。復元抽出で重複行を残さない\*/**

```
PROC SURVEYSELECT DATA=ksug.customer
  METHOD=URS
  N=2000
  OUT=sample6;
RUN;QUIT;
```

	顧客番号	年齢	性別	FID	売上金額	婚姻状況	住居種類	Number Hits (Default=1)
1	42	49	女性	1	5200	既婚	マンション	1
2	55	37	女性	1	5	既婚	社宅	1
3	76	24	女性	0	22944	既婚	マンション	1
4	90	32	女性	0	5000	既婚	マンション	1
5	100	58	女性	1	1400	既婚	アパート	1
6	107	28	女性	0	3200	未婚	マンション	1
7	110	27	女性	0	4800	未婚	マンション	1
8	114	35	女性	0	5	既婚	雑居	1
9	117	38	女性	0	1000	既婚	マンション	1
10	127	48	女性	0	5	既婚	雑居	1
11	136	29	女性	0	5	未婚	雑居	1
12	136	39	女性	0	8400	未婚	マンション	1
13	144	27	女性	0	5	未婚	マンション	1
14	150	27	女性	1	7000	既婚	マンション	1
15	152	38	女性	0	1000	既婚	マンション	1

今回のプログラムは、手法が「METHOD=URS」であることが大きな違い、URSはUnrestricted Random Samplingの事。

復元抽出を行っているため、何件かは複数回抽出されている。しかし何度も同じ標本が選ばれると、結果が歪んでしまうことも考えられるので、今回はデフォルト設定のまま複数回選択されたものは1件だけ残して削除する方法を取っている。

何件重複しているかは、新たに追加される変数「NumberHits」を見ると分かる。よって、このプログラムの結果は何行が削除されているので、2000件と指定したが実際には1917件しか抽出されていない。

## 単純無作為抽出 3

**/\*件数を指定。復元抽出で重複行を残す\*/**

```
PROC SURVEYSELECT DATA=ksug.customer
  METHOD=URS
  N=2000
  OUT=sample7
  OUTHITS;
RUN;QUIT;
```

	顧客番号	年齢	性別	FID	売上金額	婚姻状況	住居種類	Number Hits (Default=1)
120	1179	35	女性	0	5	未婚	マンション	1
121	1189	35	女性	0	5	既婚	マンション	1
122	1206	42	女性	1	5	既婚	マンション	1
123	1211	15	女性	0	5	既婚	マンション	1
124	1211	15	女性	0	5	既婚	マンション	1
125	1229	18	女性	0	5	既婚	マンション	1
126	1241	36	女性	0	5	既婚	一戸建て	1
127	1248	20	女性	0	8000	既婚	社宅	1
128	1248	20	女性	0	8000	既婚	社宅	2
129	1248	20	女性	0	8000	既婚	社宅	2
130	1253	30	女性	0	8000	既婚	マンション	1
131	1286	47	女性	0	5	既婚	雑居	1
132	1286	36	女性	0	5	既婚	学生寮	1
133	1289	30	女性	0	5	既婚	マンション	1
134	1289	30	女性	0	5	既婚	マンション	1
134	1289	30	女性	0	5	既婚	マンション	1

今回のプログラムも同じように手法が「METHOD=URS」である。ただ今回は「OUTHITS」オプションを付けている。

復元抽出で同じ標本が何度抽出されても、削除はせずにサンプルとして出力する。何件重複しているかは、新たに追加される変数「NumberHits」を見ると分かる。

よって、このプログラムの結果は重複しても削除される事が無いので、指定した2000件が出力されている。

## 単純無作為抽出 4

```
/*割合を指定。「N=」を「RATE=」に変更*/
PROC SURVEYSELECT DATA=ksug.customer
METHOD=SRS
RATE=0.1
OUT=sample8;
RUN;
QUIT;
```

単純無作為抽出をする時に、そのサンプルサイズを「N=2000」といった方法ではなく、全体の10%といった割合を指定して抽出を行う方法。

「N=」ではなくて、「RATE=0.1」と指定すると10%サンプリングを行える。

観測番号	年齢	性別	州	売上金額	婚姻状況	住居形態
1	4	女性	0	0	結婚	マンション(借)
2	16	女性	0	62700	結婚	借家
3	30	男性	0	11800	未婚	マンション(借)
4	40	女性	1	0	未婚	社宅
5	42	女性	1	6000	結婚	マンション(借)
6	40	女性	0	6800	結婚	学生寮
7	46	男性	0	0	未婚	社宅
8	48	男性	0	1400	結婚	一戸建て
9	56	女性	0	2500	未婚	借家
10	60	女性	0	19000	未婚	マンション(借)
11	64	男性	0	4800	結婚	マンション(借)
12	66	女性	0	0	未婚	アパート
13	67	女性	1	20800	未婚	借家
14	77	女性	0	0	未婚	社宅
15	81	女性	0	0	未婚	マンション(借)

17

## 系統抽出(等間隔抽出法)

```
/*件数を指定*/
PROC SURVEYSELECT DATA=ksug.customer
METHOD=SYS
N=2000
OUT=sample9;
RUN;
QUIT;
```

系統抽出を行う際には、「METHOD=SYS」と指定を行う。これはSystematic Random Samplingの略となっている。

全件で20,000件のデータから2,000件のサンプルを作るが、等間隔で抽出するので、 $200000 / 2000 = 10$ という計算を行い、10件ずつ抽出してくる。初期値はシステム時間などを利用して、決定する。

観測番号	年齢	性別	州	売上金額	婚姻状況	住居形態
1	30	女性	0	0	未婚	学生寮
2	30	男性	0	10000	結婚	社宅
3	30	女性	0	0	未婚	借家
4	40	女性	1	0	未婚	社宅
5	50	男性	0	0	結婚	アパート
6	40	女性	0	10000	未婚	マンション(借)
7	60	男性	0	0	結婚	マンション(借)
8	60	女性	0	20000	結婚	マンション(借)
9	50	女性	0	6800	結婚	学生寮
10	60	女性	1	16000	未婚	アパート
11	110	女性	0	0	結婚	マンション(借)
12	120	男性	0	4800	結婚	一戸建て
13	130	女性	0	0	未婚	マンション(借)
14	140	女性	1	14000	結婚	学生寮
15	150	女性	0	0	結婚	一戸建て

データがランダムに並んでいる、もしくは北海道～沖縄県といった形で並んでいる場合はこれを反映することができる。ただし、何か特定のルールで並んでいる場合は、初期値によって全く傾向が違ったサンプルが抽出される恐れがあるので注意が必要。

18

## 層別抽出法 1

```

/*性別別に1000件ずつ持ってくる場合。
 [男性][女性][不明]なので合計3000名*/
PROC SORT DATA=ksug.customer;
  BY gender;
RUN;
PROC SURVEYSELECT DATA=ksug.customer
  METHOD=SRS N=1000 OUT=sample10;
  STRATA gender;
RUN;QUIT;
  
```

性別				
GENDER	度数	パーセント	累積 度数	累積 パーセント
女性	1000	33.33	1000	33.33
男性	1000	33.33	2000	66.67
不明	1000	33.33	3000	100.00

上記のプログラムでは、性別別に1000件ずつ抽出し、サンプルの作成を行っている。  
 「STRATAステートメント」で変数「GENDER」を指定する事により、変数「GENDER」内の水準  
 を利用して、重み付けを行い同じ件数で抽出を行っている。  
 出力されたデータには重み付けの結果の変数などが追加されている。

またデータは「STRATAステートメント」で指定する変数でソートを行う、もしくはINDEXを作成し  
 ておく必要がある。

19

## 層別抽出法 2

```

/*性別別に10% 件ずつ持ってくる場合*/
PROC SORT DATA=ksug.customer;
  BY gender;
RUN;
PROC SURVEYSELECT DATA=ksug.customer
  METHOD=SRS
  RATE=0.1
  OUT=sample11;
  STRATA gender;
RUN;quit;
  
```

性別				
GENDER	度数	パーセント	累積 度数	累積 パーセント
女性	807	40.33	807	40.33
男性	821	41.03	1628	81.36
不明	373	18.64	2001	100.00

上記のプログラムでは、性別別に10%件ずつ抽出し、サンプルの作成を行っている。  
 結果的に20,000件のうち2001件のサンプルを抽出している。  
 この場合も前ページと同じように、出力されたデータには重み付けの結果の変数な  
 どが追加されている。またデータは「STRATAステートメント」で指定する変数でソ  
 ートを行う、もしくはINDEXを作成しておく必要がある。

20

## 繰り返してサンプリングを行うには？

**/\*100オブザベーションの10のサンプルを作成するには？\*/**

```
PROC SURVEYSELECT DATA=ksug.customer
METHOD=SRS
N=100
OUT=sample12
REP=10;
RUN;QUIT;
```

Customer Replicate Number	顧客番号	年齢	性別	F10	売上金額	増殖状況	増殖種類
36	15764	26	不明	8	2500	増殖	ランダム(増)
37	16780	24	不明	8	130700	増殖	ランダム(増)
38	16830	27	不明	8	48800	増殖	ランダム(増)
39	16862	76	不明	8	28000	増殖	ランダム(増)
180	16964	30	不明	8	8	増殖	ランダム(増)
181	485	32	女性	1	4000	増殖	ランダム(増)
182	880	28	女性	1	8400	増殖	増殖
183	1325	48	女性	8	6000	増殖	一戸建て
184	1330	37	女性	8	8	増殖	増殖
185	1381	33	女性	1	8	増殖	住宅
186	2012	38	女性	8	3400	増殖	ランダム(増)
187	2322	44	女性	8	48300	増殖	ランダム(増)
188	2489	42	女性	1	15300	増殖	住宅
189	2425	38	女性	8	12000	増殖	増殖
190	2429	32	女性	8	8	増殖	ランダム(増)

左記は母集団のデータから無作為に、100オブザベーションの10サンプルを作成する方法となっています。

プログラムとしては普通の単純無作為抽出の場合と同じですが、「REP」オプションに値を設定する事により、その回数分サンプリングを行います。

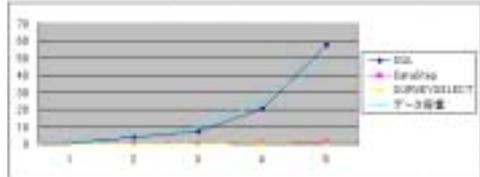
この処理結果は別々のデータセットではなく、一つのデータセットとして出力されるので、変数「Replicate」に繰り返しの回数が保存されます。

## パフォーマンスの比較

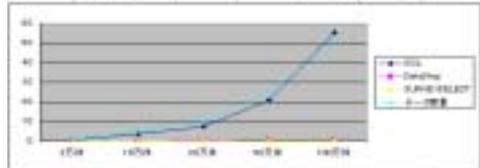
-SQL プロシジャ、Dataステップ、SURVEYSELECT プロシジャでの比較-

- ◆ テスト結果に関して
  - テストパターンは2通りを実施した。
  - 1つ目は2万~100万件まで母集団を増やしていき、その中で2000件を抽出する時間をSQLプロシジャ、Dataステップ、SURVEYSELECTプロシジャで比較した。
  - ふたつ目は上記とほぼ同じであるが、サンプルサイズを母集団の10%と指定を行った。
- ◆ データステップとSURVEYSELECTプロシジャのスピードが目立つが、記述の簡単さ、手法の多さを考えるとSURVEYSELECTプロシジャは便利に利用する事ができる。

抽出件数	SQL	Dataステップ	SQLプロシジャ	Dataステップ	SQLプロシジャ	Dataステップ
2万	0.07	0.1	0.18	0.28	0.05	0.05
10万	0.04	0.12	0.22	0.29	0.06	0.06
20万	0.05	0.1	0.17	0.27	0.07	0.07
母集団の10%	1	8	1.6	2.8	1.0	1.0



抽出件数	SQL	Dataステップ	SQLプロシジャ	Dataステップ	SQLプロシジャ	Dataステップ
2万	0.04	0.05	0.08	0.07	0.04	0.04
10万	0.04	0.12	0.15	0.20	0.05	0.05
20万	0.04	0.14	0.21	0.28	0.06	0.06
母集団の10%	1	8	1.6	2.8	1.0	1.0



## SASによる主成分分析について

## 頂いた質問の背景

- ◆ ユーザーの方はEnterpriseMinerを利用しており回帰分析等の手法を使っている。
- ◆ データには数値項目の変数を多く保持しており、回帰分析を行う時に変数間の共線性の問題がある。それぞれの変数を幾つかのまとめて、その結果を利用して回帰分析を行いたいということであった(主成分回帰)。
- ◆ その際に主成分分析の結果、主成分スコアの変数で有効なもの(固有値 $\geq 1$ )だけを自動的に残すようなマクロがあると便利なんだけど……。
- ◆ ということでサンプルを作成をしました。
- ◆ ちなみに今回のデータと質問を頂いた会社とは全く関係ありません。

## 主成分分析とは？

- ◆ 主成分分析とは？
  - 主成分分析とは、幾つかの変数から主成分というお互いに無相関な合成変数を作るための手法。これにより複数の特徴を表す総合的な指標を作成する事が出来る。
- ◆ 特徴は？
  - 生成された主成分スコアはお互いに無相関である。
  - 主成分スコアは平均0、分散1に標準化されている。
  - 生成された主成分スコアは複数生成されるが、いくつ利用するかは以下のような判断基準がある。解釈を考えると3~7つが限度。(目安程度とお考えを・・・)
    - ◆ 固有値を目安にする方法・・・基準は1以上
    - ◆ 累積寄与率を目安とする方法・・・基準は70%~80%程度
    - ◆ スクリーンプロットを目安にする方法・・・固有値をパレート図などを用いて大きい順に並べ、急激に落ち込む手前の主成分。
  - 主成分をいかにして解釈するかが重要。(意見は色々あると思いますが・・・)
  - 他の手法と同様だが、色々なグラフや手法と組み合わせると有効な結果が出てくると考えられる。

## 今回利用するデータは？

№	変数	ラベル
1	COMPANY	会社名
2	SUNRI	業種大分類
3	GYO	業種中分類
4	SISAN_RYUDO	流動資産合計
5	SISAN_TANAGROSHI	棚卸資産合計
6	SISAN_YUKE(KOTEI)	有形固定資産合計
7	SISAN_MUKE(KOTEI)	無形固定資産合計
8	SISAN_TOSHI_SONOTA	投資・その他資産合計
9	FUSAI_RYUDO	流動負債合計
10	FUSAI_KOTEI	固定負債合計
11	SHON	資本合計
12	SAIAGE	売上高
13	KOKOKU_SENDOENJI	広告費・宣伝費
14	HANKANHI	販売費及び一般管理費
15	RIEKI_EIGYO	営業利益
16	EIGYOGAI_SYUEKI	営業外収益合計
17	EIGYOGAI_HIYO	営業外費用合計
18	RIEKI_KEIJO	経常利益
19	TOKUBETSU_RIEKI	特別利益合計
20	TOKUBETSU_SONSHITSU	特別損失合計
21	HAITOKIN	年間支払配当金
22	RIEKI_TOXI	当期利益

- ◆ 左図のような変数の構成をしているデータを使って、一番下の当期利益を予測できるようなモデルを作成したい。
- ◆ 一般的に当期利益を算出するためのロジックは決まっているため、ある一定の勘定項目が分かれば算出できてしまう。また金額といった値だけではなく、「流動比率」「売上高対総資産回転率」等の財務指標や格付けのデータなどがあれば面白い。
- ◆ 多くの数値変数があり、かつ微妙に絡み合っているが、業種によっては大きく構成が違ってくる。
- ◆ この辺りを主成分分析を利用して解析するサンプルをご紹介し、質問に出てきた主成分回帰を行えるデータセットを作成する。

## 主成分分析の実行と結果 1

```
proc princomp data=KSUG.ENTERPRISE
  out=OUT_PRIN;
  var SISAN_RYUDO--HAITOKIN;
run;quit;
```

Eigenvalues of the Correlation Matrix			
	Eigenvalue	Difference	Proportion
1	8.3881952	5.6498195	1.9377
2	2.3945581	4.0218872	0.7492
3	1.5127254	4.4781514	0.2886
4	1.3948285	4.4350782	0.0604
5	0.8935587	4.0529899	0.0412
6	0.8081531	4.0310795	0.0279
7	0.4258275	4.1818875	0.0125
8	0.3141981	4.1388081	0.0074
9	0.2191486	4.0419442	0.0052
10	0.1671882	4.0189143	0.0034
11	0.1278014	4.0406741	0.0019
12	0.1071042	4.0424023	0.0014
13	0.1448585	4.0218124	0.0009
14	0.1211511	4.0496701	0.0007
15	0.1151035	4.0158581	0.0005
16	0.1001055		0.0003

データを利用して主成分分析を行う。今回は目的変数以外の変数を利用。出力結果としては、変数の基礎統計量と相関行列、固有値、固有ベクトルなどが出力される。固有値(Eigenvalue)の出力結果を表示している。

先ほどの主成分の残差基準である、固有値が1以上、累積寄与率(Cumulative)が80%と考えると第3主成分、もしくは第4主成分まで利用すればデータの説明はできそう。

「OUT=」オプションで主成分スコアのついたデータを書き出している。

## 主成分分析の実行と結果 2

- 前ページのプログラムにより、アウトプットとして入力データセットに主成分スコアが付けられたデータセットが作成される。この主成分スコアの変数には、プログラム中「var」オプションで指定された分だけ作成される。
- 今回は16変数追加されている。データ件数が多かったり、もっと変数が多い場合などはもう一度プログラムを作って、前ページのアウトプット画面を見ながら、DROPする変数を決めてデータセットを作り直すのも面倒となってくる。

データセット  
「OUT\_PRIN」

Prin1 ~ Prin16  
まで元のデータ  
セットに追加され  
て出力

## データ可能のためのマクロ紹介

```

%macro prin_dat1(indata,outdata,input,eigen);
ods trace on;
ods output Eigenvalues=Eigenvalues;
proc princomp data=&indata out=&outdata;
var &input;
run;quit;
ods output close;
data Eigenvalues;
set Eigenvalues;
if eigenvalue < &eigen then do ;
n+1;
call symput('drop_prin' || left(put(n,best.)),
compress(left('prin'|number))););
call symput('n_data',n);
end;
run;
data &outdata;
set &outdata(drop = %do i=1 %to &n_data; &&drop_prin&i %end;); ;
run;
%mend;
%prin_dat1(KSUG.ENTERPRISE,OUT_PRIN, S1SAN_RYUDO--HAITOKIN,1);

```

- ここで説明を行ってきた事を行おうとした場合、有効な主成分を残すためには、出力結果を見ながら、その後に関係分析時に主成分を選択する必要があるので、また不必要な変数もあることからデータ容量が大きくなってしまふ。
- そこで指定した固有値以上の変数だけを残すようにする。
- こうするとPRINCOMPプロシジャのOUTオプションに指定されたデータセットには必要な主成分スコアと元にあった変数だけが出力される。
- マクロ全体の流れとしては、可変となるインプットのデータセット、出力のデータセット、分析変数、固有値の閾値を設定すれば自動的に実行されるようになっている。
- プログラムの ~ に関して説明を加えていく。

29

## 主成分分析のアウトプットとデータセット化

```

ods trace on;
ods output Eigenvalues=Eigenvalues;
proc princomp data=&indata out=&outdata;
var &input;
run;quit;
ods output close;

```

名前: Eigenvalues  
ラベル: Eigenvalues of the Correlation Matrix  
テンプレート: Stat.PrinComp.Eigenvalues  
データラベル: Eigenvalues  
パス: Princomp.Eigenvalues

Number	Eigenvalue	Orthogonal	Principal	Component	Component
1	0.22888933	0.64319966	0.0221	0.8224	
2	0.20488591	0.82719977	0.1477	0.8527	
3	0.15217684	0.47819276	0.0882	0.7684	
4	0.08847082	0.47819276	0.0594	0.8788	
5	0.08847082	0.47819276	0.0594	0.8788	
6	0.08847082	0.47819276	0.0594	0.8788	
7	0.08847082	0.47819276	0.0594	0.8788	
8	0.08847082	0.47819276	0.0594	0.8788	
9	0.08847082	0.47819276	0.0594	0.8788	
10	0.08847082	0.47819276	0.0594	0.8788	
11	0.08847082	0.47819276	0.0594	0.8788	
12	0.08847082	0.47819276	0.0594	0.8788	
13	0.08847082	0.47819276	0.0594	0.8788	
14	0.08847082	0.47819276	0.0594	0.8788	
15	0.08847082	0.47819276	0.0594	0.8788	
16	0.08847082	0.47819276	0.0594	0.8788	

PRINCOMPプロシジャなどの統計量を活用する場合は、「OUT」オプション、「OUTSTAT」オプションなどが利用できるが、今回は固有値だけのデータが欲しいのでODSの機能を利用する。

ods trace on;  
PRINCOMPプロシジャでどのようなアウトプットが出せて、そのアウトプットをデータセット化するための名前を知る必要がある。その場合には「ods trace on;」と記述する事により、ログにどのような統計量が出力されるか結果が分かる。

ods output Eigenvalues=Eigenvalues;  
今回は固有値のデータセットを作成したいので、「ods output Eigenvalues」という部分で、固有値をEigenvaluesというデータセットに出力する指定をしている。

出力されるデータセットを見てみると、主成分の固有値で1未満の変数は「Prin5」以降だということになるので、自動的にこの変数をセットして削除することを行う。

30

## 削除する変数の決定と変数へのセット

```
data _NULL_;
set Eigenvalues;
if eigenvalue < &eigen then do ;
n+1;
call symput('drop_prin' || left(put(n,best.)),
compress(left('prin'||number)));
call symput('n_data',n);
end;
run;
```

マクロ変数名	値
drop_prin1	Prin5
drop_prin2	Prin6
drop_prin3	Prin7
drop_prin4	Prin8
drop_prin5	Prin9
drop_prin6	Prin10
drop_prin7	Prin11
drop_prin8	Prin12
drop_prin9	Prin13
drop_prin10	Prin14
drop_prin11	Prin15
drop_prin12	Prin16

前ページのデータセット「Eigenvalues」を読み込んで、固有値が指定した値以下の物だけマクロ変数に格納をしていき、その件数もセットしておく。今回は固有値を1にセットしているので第5主成分以下が格納されている。

マクロ変数には右の表のような形で格納されている。

マクロ変数名	値
n_data	12

## 余分な主成分スコア変数の削除

```
data &outdata;
set &outdata(drop =
%do i=1 %to &n_data;
&&drop_prin&i;
%end;);
run;
```

マクロ変数名	値
n_data	12

マクロ変数名	値
drop_prin1	Prin5
drop_prin2	Prin6
drop_prin3	Prin7
drop_prin4	Prin8
drop_prin5	Prin9
drop_prin6	Prin10
drop_prin7	Prin11
drop_prin8	Prin12
drop_prin9	Prin13
drop_prin10	Prin14
drop_prin11	Prin15
drop_prin12	Prin16

変数名	変数値	変数名	変数値	変数名	変数値
1	517	200	-0.29175821	31002	0.00000000
2	520	199	-0.19328173	31003	0.00000000
3	294	198	-0.19328173	31004	0.00000000
4	294	198	-0.18447027	31005	0.00000000
5	180	400	-0.19328173	31006	0.00000000
6	180	400	-0.18447027	31007	0.00000000
7	180	400	-0.18447027	31008	0.00000000
8	180	400	-0.18447027	31009	0.00000000
9	180	400	-0.18447027	31010	0.00000000
10	180	400	-0.18447027	31011	0.00000000
11	180	400	-0.18447027	31012	0.00000000
12	180	400	-0.18447027	31013	0.00000000
13	180	400	-0.18447027	31014	0.00000000
14	180	400	-0.18447027	31015	0.00000000
15	180	400	-0.18447027	31016	0.00000000

最後にPRINCOMPプロシジャで出力され、まだ多くの主成分スコア変数が残されているデータを読み込む。

マクロに格納されている固有値が1未満の変数を呼出して、「SET」ステートメントの「DROP」オプションにセットをしていき、いらない変数を削除する。そして最終的なデータセットが書き出される。

これにより必要な変数だけが残されたデータセットが完成する。

## 他には何か出来ないか？

- ◆ 今回はサンプリングと主成分分析を題材にしたプロシジャの結果のデータ二次加工についてご紹介しました。
- ◆ 一般的に企業で使われるデータの分布は非常歪んでおり、このような場合にモデルの妥当性を評価する必要があり、有名な方法としてジャックナイフ法、ブートストラップ法などがある。
- ◆ ふたつのテーマを組み合わせる事により、データからサンプリングを行い、主成分分析を繰り返して行って、結果の安定性やバラツキを知る事ができるような処理も組み込めると考えられる。



The Power to Know™