

# 統計学第14/15講

明治大学情報コミュニケーション学部

後藤 晶

akiragoto@meiji.ac.jp

## 今日のお話

前回の復習

2要因分散分析モデル(交互作用あり)

2要因分散分析モデル(交互作用なし)

## 前回の復習

## 分散分析とは

分散分析とは、「3群以上の分散に差があるかどうか」を比較・分析するための方法です。その後「多重比較」という手法を用いて、「3群以上の平均値の差があるかどうか」を明らかにします。この授業では「1元配置分散分析」および「2元配置分散分析」というものについて説明します。いずれについても、説明変数が因子データ、応答変数が数値データとなります。

- ▶ 1元配置分散分析：「地域によって、主観的幸福度の分散・平均値が異なる」などのような、1つの要因によって影響を受けるかどうかを分析する手法です。
- ▶ 2元配置分散分析：「地域と未婚・既婚によって分散・平均値が主観的幸福度が異なる」、「地域と子の有無によって主観的幸福度が異なる」などのような、2つの要因によって影響を受けるかどうかを分析する手法です。

分散分析を一般線形モデルの枠組みで説明すると、平均値の推定がベースとなりますが、以下のように理解することができます。ここでは、「3つの群の影響を受ける」場合について、モデル式を元に説明します。また、以下では「分散分析モデル」という表現をします。

- ▶ 個人的には一般線形モデルの枠組みの方が理解しやすいと思っています。

$$Y_i = \beta_1 X_1 + \beta_2 X_2 + \alpha + \epsilon_i$$

このモデルでは $X_1$ と $X_2$ はそれぞれ(1, 0)の値を取る「ダミー変数」です。しかし、これでは $\beta$ が2つしかありません。しかし、これだけで3つの群を表すことができます。以下には3つの条件についてモデル式を書き入れてあげたいと思います。

▶  $X_1 = 1$  と  $X_2 = 0$  の場合

$$Y_i = \beta_1 + \alpha + \epsilon_i$$

- この場合, ある因子  $X_1$  によって, 傾きが変化することを示しています.

▶  $X_1 = 0$  と  $X_2 = 1$  の場合

$$Y_i = \beta_2 + \alpha + \epsilon_i$$

- この場合, ある因子  $X_2$  によって, 傾きが変化することを示しています.

▶  $X_1 = 0$  と  $X_2 = 0$  の場合

$$Y_i = \alpha + \epsilon_i$$

このモデルについて，平均値が異なるかどうかを調べます．特に，分散分析の場合は「分散分析表」と呼ばれるものを出して評価してあげます．

## 分散分析モデルの例

- ▶ テストの点数がクラスによって異なる.

$$Y_i = \beta_1 X_1 + \beta_2 X_2 + \alpha + \epsilon_i$$

- ▶  $X_1 = 1$ と $X_2 = 0$  : Bクラス

$$Y_i = \beta_1 + \alpha + \epsilon_i$$

- ▶  $X_1 = 0$ と $X_2 = 1$  : Cクラス

$$Y_i = \beta_2 + \alpha + \epsilon_i$$

- ▶  $X_1 = 0$ と $X_2 = 0$  : Aクラス

$$Y_i = \alpha + \epsilon_i$$



## 仮説を立てる

さて、それでは仮説を立ててみましょう。今回分析するテーマは「主観的幸福度(SUB\_HAP)が地域(SUB\_ARE)によって異なる」かどうかを分析します。一要因分散分析の場合は以下のような仮説を立てます。

- ▶ 対立仮説：主観的幸福度の平均値は地域によって異なる。
- ▶ 帰無仮説：主観的幸福度の平均値は地域によって異なるとは言えない。

この2つの仮説のもとに分析を行ないます。

## 分析のモデル式

今回の分析には、以下のモデルを前提とします.

$$(SH) = \beta_1(Hokkaido\_dum) + \beta_2(Tohoku\_dum) + \beta_3(Chubu\_dum) + \beta_4(Kinki\_dum) + \beta_5(Chugoku\_dum) + \beta_6(Shikoku\_dum) + \beta_7(Kyushu\_dum) + \alpha + \epsilon_i$$

- ▶ なお、このモデルではそれぞれの値は1か0の値しか取りません.
- ▶ ex.東北地方のデータである場合には、東北ダミーが1、それ以外のダミー変数は0を取ります.
- ▶ また、すべてのダミー変数が0の場合はコントロール群となる関東地方の値を示しています.

## 平均値をプロットする

さて，例によってggplotguiを使いましょう．

以下のコードはConsole（コンソール）に直接打ち込みます．

```
library(ggplotgui)  
ggplot_shiny()
```

そうすると新しいウィンドウが開きます．

以下の通りの作業をしましょう.

- ▶ **"Data upload"**をクリック
- ▶ datasetをコピーする
- ▶ **"Paste Data"**にペーストをする
- ▶ ggplotタブへ
- ▶ **"Type of graph:"**は**"Dot + Error"**, Y-variableは**"SUB\_HAP"**, X-variableは**"ARE"**を設定
- ▶ **"Confidence Interval:"**を95%にする.
- ▶ R-codeタブへ行って, 以下のコードのうち, 真ん中のみを以下にする. -また, コード内の**df**を**exdataset**に変える.

▶ こんな感じのコードができます.

```
# You need the following package(s):
```

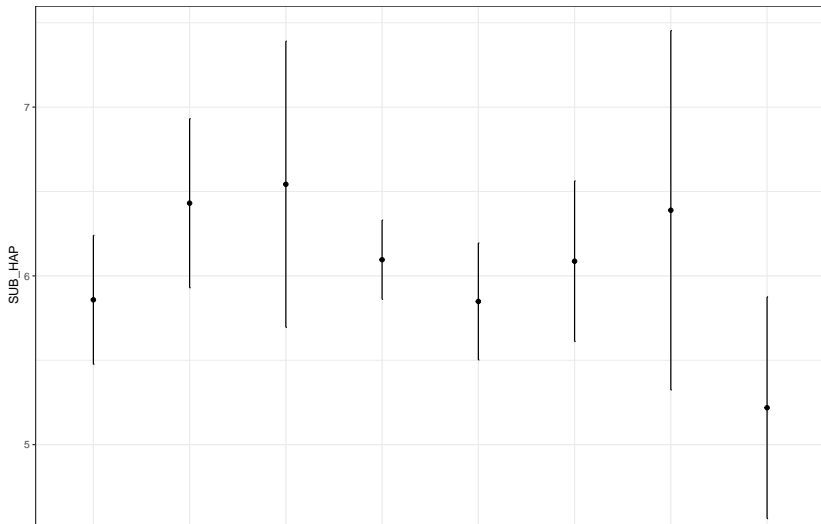
```
library("ggplot2")
```

```
# The code below will generate the graph:
```

```
graph <- ggplot(exdataset, aes(x = ARE, y = SUB_HAP)) +  
  geom_point(stat = 'summary', fun.y = 'mean') +  
  geom_errorbar(stat = 'summary', fun.data = 'mean_se',  
               width=0, fun.args = list(mult = 1.96)) +  
  theme_bw()
```

そうすると, こんなグラフが算出されます.

graph



このグラフを見る限り，地域ごとに差があるかどうかはわかりません．以前，平均値を算出してみたことがありましたが，今回はそれぞれが「統計的に差がある」ということが言えるかどうかを考えたいと思います．

## 分析をやってみる

さて，分散分析モデルを作成してみましょう．

```
"arehap_model"  
arehap_model<-lm(SUB_HAP ~ ARE, data = exdataset)
```



```
#
summary(arehap_model)

##
## Call:
## lm(formula = SUB_HAP ~ ARE, data = exdataset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.5429 -1.4308  0.1515  1.9043  4.7813
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  5.858108   0.192158  30.486  <2e-16 ***
## AREChugoku   0.572661   0.347849   1.646   0.1000
## AREHokkaido  0.684749   0.439390   1.558   0.1195
## AREKanto     0.237637   0.226845   1.048   0.2951
## AREKinki     -0.009623   0.264660  -0.036   0.9710
```

- ▶ 出力結果が入り切らないのでCoefficientsだけ示します.

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	6.095745	0.120558	50.563	< 2e-16	***
AREHokkaido	0.447112	0.413125	1.082	0.27941	
ARETohoku	-0.876995	0.316105	-2.774	0.00564	**
AREChubu	-0.237637	0.226845	-1.048	0.29510	
AREKinki	-0.247260	0.218299	-1.133	0.25764	
AREChugoku	0.335025	0.314020	1.067	0.28629	
AREShikoku	0.293144	0.564036	0.520	0.60338	
AREKyushu	-0.008788	0.271909	-0.032	0.97422	

- ▶  $\alpha$ は6.095745である.
- ▶ 関東地方と比べて、東北地方の主観的幸福度が低い.
  - 実は昔から言われている結果.
  - 東日本大震災の影響? という声もあったが逆で、東日本大震災によって幸福度が改善したとも言われている.

## 分析結果の解釈

- ▶ さらに，モデル式による分析結果を出力しました．この結果が示しているのは以下のようなことです．

$$\begin{aligned}(SH) = & 0.447112 * (Hokkaido\_dum) - 0.876995 * (Tohoku\_dum) \\ & 0.237637 * (Chubu\_dum) - 0.247260 * (Kinki\_dum) + \\ & 0.335025 * (Chugoku\_dum) + 0.293144 * (Shikoku\_dum) - \\ & 0.008788 * (Kyushu\_dum) + 6.095745 + \epsilon_i\end{aligned}$$

- ▶ 今度はモデル式についても同じよう出力してあげましょう．
- ▶ 回帰分析やt検定と同じです．

## 分散分析表の出力

#

```
anova(arehap_model)
```

```
## Analysis of Variance Table
```

```
##
```

```
## Response: SUB_HAP
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
```

```
## ARE           7    75.1  10.7238   1.9623 0.05729 .
```

```
## Residuals 955 5218.9   5.4648
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
```

## 分散分析表の読み方：

### Analysis of Variance Table

- ▶ 分散分析表です．分散分析の結果を示しています．

Response: SUB\_HAP

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
ARE	7	75.1	10.7238	1.9623	0.05729 .
Residuals	955	5218.9	5.4648		

- ▶ Dfは自由度を示しています.
- ▶ Sum Sqは平方和
- ▶ Mean Sqは平均平方
- ▶ F valueはF値
- ▶ Pr(>F)はp値を示しています.

Response: SUB\_HAP

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
ARE	7	75.1	10.7238	1.9623	0.05729
Residuals	955	5218.9	5.4648		

- ▶ 応答変数はSUB\_HAPです.
- ▶ AREの自由度（分子自由度）は7：全部で8地域ある→N-1が自由度
  - モデル式の $\beta$ （パラメータ）の数と一致している.
  - DFはDegree of Freedom
- ▶ AREのF値は1.9623, P値は0.05729
- ▶ Residualsの自由度（分母自由度）は955：全部で963個のデータがあり, モデル式の $\beta$ （パラメータ）で7つ, さらにもう1地域（= $\alpha$ で使われる）を引いたもの.

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

- ▶ p値の大きさを示す記号.
- ▶ 0 -0.001では\*\*\*で表される.
- ▶ 0.001-0.01では\*\*で表される.
- ▶ 0.01-0.05では\*で表される.
- ▶ 0.05-0.1では.で表される.
- ▶ 0.1-1では何にもありません.



## 書き方

- ▶ 主観的幸福度は地域によって異なるかを分析した。その結果、 $F(7, 955)=1.9623(p<.10)$ であり、10%水準で有意であることが示されている。したがって、主観的幸福度は居住地域によって異なる傾向にあることが示されている。
  - 分散分析表を合わせて示してあげましょう。
  - ちなみに、心理学などでは有意水準を5%に設定されることが多い。
  - 経済学系では10%水準を採用することもある。
  - いずれにしろ、分析の前に有意水準を設定する必要がある。

## 結果をきれいに表記しよう.

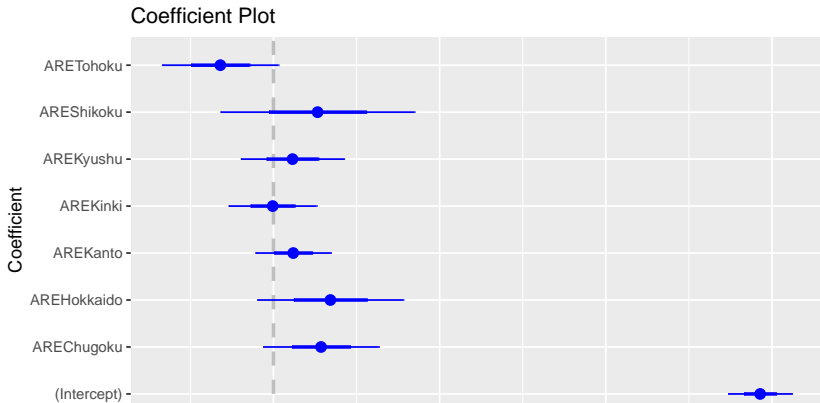
- ▶ パッケージhuxtableの中にhuxregという関数があります.

```
library(huxtable)  
huxreg(arehap_model)
```

## 結果をきれいに表記しよう.

- ▶ パッケージcoefplotを使って各係数の大きさをグラフで示す.

```
library(coefplot)
coefplot(arehap_model)
```



## 結果をきれいに表記しよう.

- ▶ パッケージstargazerの中にあるstargazerという関数を使うとxls形式で出力できます.

```
library(stargazer)
stargazer(arehap_model, type = "html", align=TRUE, title =
  "", out = "arehap_model.xls")
```

## 自由度とは

- ▶ 自由度 =  $n - p$ 
  - $n$ : 標本の大きさ
  - $p$ : 推定されたパラメータの数
- ▶ 自由度 =  $n - q - 1$ 
  - $n$ : 標本の大きさ
  - $q$ : モデル式で推定されたパラメータ ( $\beta$ ) の数
  - 1は ( $\alpha$ ) の分

## 要約

- ▶ 一般線形モデルによる分散分析モデル
  - ダミー変数が複数あるような状況を前提とする.

```

<-lm(      ~      ,
          data =      )

      t
summary(    )

anova(    )

```

## 2要因分散分析モデル(交互作用あり)

## 2要因分散分析(交互作用あり)

- ▶ 続いて、2要因分散分析に進みたいと思います。2要因分散分析とは、複数の要因による影響を分析するものです。例えば、主観的幸福度は子の有無(1,0のダミー変数)だけでなく、結婚しているか否か(1,0のダミー変数)によっても影響を受ける可能性があります。これを用いると「子がいない未婚者」「子がいない既婚者」「子がいる未婚者」「子がいる既婚者」の計4つの状態があります。
- ▶ したがって、これらが影響を与えているかどうかを明らかにするために、いずれの要因についても投入したモデル式について考えたいと思います。ここでは、次のようなモデル式を考えたいと思います。

$$Y_i = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 * X_2 + \alpha + \epsilon_i$$

このモデル式によって、「4つの状態」を分析することができ



▶  $X_1 = 1$ と $X_2 = 0$ の場合

$$Y_i = \beta_1 * 1 + \beta_2 * 0 + \beta_3 * 1 * 0 + \alpha + \epsilon_i$$

$$Y_i = \beta_1 + \alpha + \epsilon_i$$

- この場合, ある因子 $X_1$ によって, 傾きが変化することを示しています.
- ex. 子がない独身者よりも, 子がいる独身者の方が幸せとか

▶  $X_1 = 0$ と $X_2 = 1$ の場合

$$Y_i = \beta_1 * 0 + \beta_2 * 1 + \beta_3 * 0 * 1 + \alpha + \epsilon_i$$

$$Y_i = \beta_2 + \alpha + \epsilon_i$$

- この場合, ある因子 $X_2$ によって, 傾きが変化することを示しています.
- ex. 子がない未婚者よりも, 子がない既婚者の方が幸せとか

►  $X_1 = 1$ と $X_2 = 1$ )の場合

$$Y_i = \beta_1 * 1 + \beta_2 * 1 + \beta_3 * 1 * 1 + \alpha + \epsilon_i$$

$$Y_i = \beta_1 + \beta_2 + \beta_3 + \alpha + \epsilon_i$$

- この場合,  $X_1$  と  $X_2$  が影響する場合の値を示していることとなります. 特に,  $X_1 * X_2$  の係数が有意になる場合は単純に  $X_1$  と  $X_2$  が同じように影響を与えているだけでなく, 組み合わせることによって効果が強まることを示しています.
- 「組み合わせることにより効果が変わる」ことを「交互作用」といいます.
- ex. 子がない未婚者よりも, 子がいる既婚者の方が幸せ
- 子どもがいることによる幸福度の改善と, 結婚していることによる幸福度の改善から予想できないくらいググッと幸せ

- ▶  $X_1 = 0$ と $X_2 = 0$ の場合

$$Y_i = \beta_1 * 0 + \beta_2 * 0 + \beta_3 * 0 * 0 + \alpha + \epsilon_i$$

\$\$ Y\_i = \alpha + \epsilon\_i \$\$

- ▶ この場合、全ての要因が影響しない場合（何らかの基準となる点）の値を示していることになります。  
— ex. 子がない未婚者の幸福度の推定値

## 仮説を立てる

さて、それでは仮説を立ててみましょう。今回分析するテーマは「主観的幸福度(SUB\_HAP)が子の有無(CHI)と結婚(MAR)によって異なる」かどうかを分析します。二要因分散分析（交互作用有り）の場合は以下のような仮説を立てます。

- ▶ 対立仮説：主観的幸福度の平均値は結婚かつ子の有無によって異なる。
- ▶ 帰無仮説：主観的幸福度の平均値は結婚かつ子の有無によって異なるとは言えない。

この仮説のもとに分析を行ないます。

## 平均値をプロットする

さて、最初のお約束です。平均値をプロットしましょう。まずは各自でやってみましょう。

さて、例によってggplotguiを使いましょう。

以下のコードはConsole（コンソール）に直接打ち込みます。

```
library(ggplotgui)  
ggplot_shiny(exdataset)
```

そうすると新しいウィンドウが開きます。

以下の通りの作業をしましょう.

- ▶ ggplotタブへ
- ▶ "Type of graph:"は"Dot + Error", Y-variableは"SUB\_HAP", X-variableは"MAR"を設定
- ▶ "Group(or colour)"をCHIに変更
- ▶ "Confidence Interval:"を95%にする.
- ▶ R-codeタブへ行って, 以下のコードのうち, 真ん中のみを以下にする. -また, コード内のdfをexdatasetに変える.

*# You need the following package(s):*

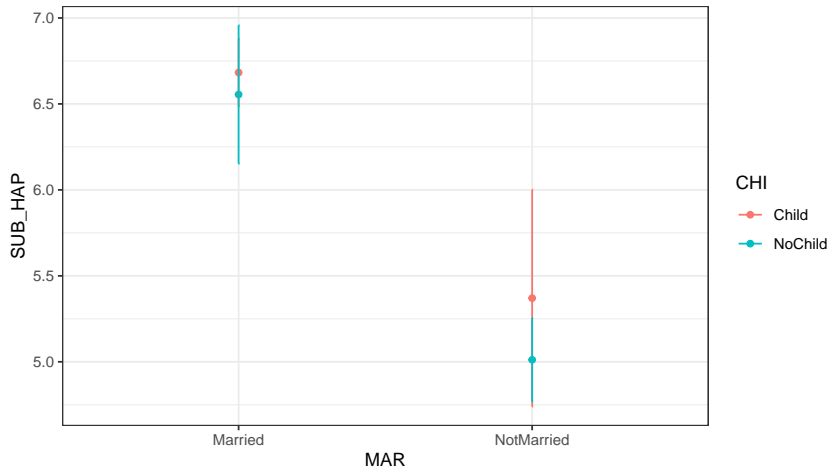
```
library("ggplot2")
```

*# The code below will generate the graph:*

```
graph <- ggplot(exdataset, aes(x = MAR, y = SUB_HAP, colour =  
  geom_point(stat = 'summary', fun.y = 'mean') +  
  geom_errorbar(stat = 'summary', fun.data = 'mean_se',  
               width=0, fun.args = list(mult = 1.96)) +  
  theme_bw()
```



graph



このグラフを見る限り，未婚者に比べて既婚者の方が主観的幸福度が高そうですが，子の有無の影響はありそうな気がしますし，なさそうな気がしますし何とも言えません．したがって，この点についても統計的に差があるのかどうかを明らかにしましょう．

## 2要因分散分析（交互作用あり）のモデル式

```
marchihap_model <- lm(SUB_HAP ~ MAR*CHI, data = exdataset)
#   MARCHIHAP_model
```

```
#
summary(marchihap_model)

##
## Call:
## lm(formula = SUB_HAP ~ MAR * CHI, data = exdataset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.6825 -1.6825  0.3175  1.3175  4.9882
##
## Coefficients:
##                                Estimate Std. Error t value Pr
t|)
## (Intercept)                6.6825     0.1054   63.419  <
## MARNotMarried             -1.3122     0.3190  -4.113  4.2
## CHINoChild                -0.1279     0.2222  -0.575
## MARNotMarried:CHINoChild  -0.2308     0.3930  -0.587
```

## 結果の書き方

この分散分析表の結果より以下のように結果を導き出すことが出来ます。 交互作用のある分散分析により、主観的幸福度は結婚および子の有無によって異なるかを分析した。その結果、結婚については $F(1, 959)=120.63(p< .001)$ であり、結婚が主観的幸福度に対して有意に影響を与えていることが明らかとなった。一方、子の有無については $F(1, 959)=1.2102(p> .05)$ 、結婚と子の有無の交互作用については $F(1, 959)=0.3448(p> .05)$ であり、有意差は認められなかった。

## 結果の解釈

この結果は以下のように解釈することができます.

$$(SH) = 1.543(Married\_dum) + 0.359(Child\_dum) - \\ 0.231(Married\_dum * Child\_dum) + 5.012$$

ただし, 以下のように変数を割り振っています.

- ▶ 結婚: 未婚→0, 既婚→1
- ▶ 子ども: 子なし→0, 子あり→1

したがって, 「未婚者かつ子なし」「未婚者かつ子あり」「既婚者かつ子なし」「既婚者かつ子あり」という4つのありえる状態について, 次のように主観的幸福度を推定することができます.

## 結果の解釈

- ▶ 「未婚者かつ子なし」

$$(SH) = 1.543 \times 0 + 0.359 \times 0 - 0.231(0 \times 0) + 5.012$$

$$(SH) = 5.012$$

- ▶ 「未婚者かつ子あり」

$$(SH) = 1.543 \times 0 + 0.359 \times 1 - 0.231(0 \times 1) + 5.012$$

$$(SH) = 0.359 + 5.012 = 5.371$$

- ▶ 「既婚者かつ子なし」

$$(SH) = 1.543 \times 1 + 0.359 \times 0 - 0.231(1 \times 0) + 5.012$$

$$(SH) = 1.543 + 5.012 = 6.555$$

- ▶ 「既婚者かつ子あり」

$$(SH) = 1.543 \times 1 + 0.359 \times 1 - 0.231(1 \times 1) + 5.012$$

$$(SH) = 1.543 + 0.359 - 0.231 + 5.012 = 6.683$$

ここから、未婚者に比べて既婚者の主観的幸福度が高いことはわかりますが、子の有無は主観的幸福度に対して影響をどうも与えなそうです。



## 分散分析（一般線形モデルによる分散分析モデルによる分析）

- ▶ 一般線形モデルによる分散分析モデル
  - ダミー変数が複数あるような状況を前提とする.
- ▶ 交互作用ありモデル：
  - 組み合わせによってパワーアップorパワーダウン. . .  
オブジェクト<-lm(応答変数 ~ 説明変数1 \* 説明変数2,  
data = データセットの名前) これについて, 回帰分析/t  
検定の際は以下のコードを使っています. summary(オブ  
ジェクト) これについて, 分散分析の際は以下のコードを  
使っています. anova(オブジェクト)

## 2要因分散分析モデル(交互作用なし)

## 2要因分散分析(交互作用なし)

今までの例題，分散分析表からは「結婚」が主観的幸福度に影響を与えることは明らかになりましたが，「子の有無」や「結婚と子の有無の交互作用」は認められませんでした。したがって，結婚をしているかどうかで主観的幸福度が高くなることは明らかとなりましたが，子がいるかどうかで主観的幸福度に影響を与えるとはいえないこと，さらに結婚しているかどうか，かつ子がいるかどうかという両者の影響が組み合わさっても影響がないことが明らかとなりました。

この結果は「未婚者かつ子なし」「未婚者かつ子あり」「既婚者かつ子なし」「既婚者かつ子あり」という4つの状態がありました。

「未婚者」に比べて、「既婚者」の主観的幸福度が高いことがわかりましたが、子の有無が与える影響と、「結婚していることかつ子の有無が与える影響」についてはあるとは言えない結果が得られました。先程の「交互作用」は「結婚していることかつ子の有無が与える影響」を示しています。

しかし、この「交互作用」が認められなかった場合は「結婚が影響しているのか？」「子の有無が影響しているのか？」のみを検討する必要があります。すなわち、「交互作用」がない場合についても検討する必要があります。そのために、「交互作用なし」の分散分析をする必要があります。

ただし，いきなり「交互作用なし」の分析，すなわち「結婚していることかつ子の有無が与える影響」はないものとして検討することもあります．これについては研究領域の違いがあるので，その領域の慣習に従ってください．

言い換えると，交互作用なしの分析では「結婚していることかつ子の有無が与える影響」という組み合わせによる特別な影響はないことを前提とした分析ということになります．

モデル式で考えると, こんな感じです.

$$Y_i = \beta_1 X_1 + \beta_2 X_2 + \alpha + \epsilon_i$$

このモデル式によって以下の4つの状態を考えることができます.

▶  $X_1 = 1, X_2 = 0$ の場合

$$Y_i = \beta_1 * 1 + \beta_2 * 0 + \alpha + \epsilon_i$$

$$Y_i = \beta_1 + \alpha + \epsilon_i$$

- この場合, ある因子 $X_1$ によって, 傾きが変化することを示しています.
- ex.既婚で, 子どもがいない人の幸福度がわかる.

▶  $X_1 = 0, X_2 = 1$ の場合

$$Y_i = \beta_1 * 0 + \beta_2 * 1 + \alpha + \epsilon_i$$

$$Y_i = \beta_2 + \alpha + \epsilon_i$$

- この場合, ある因子 $X_2$ によって, 傾きが変化することを示しています.
- ex.未婚で, 子どもがいる人の幸福度がわかる.



▶  $X_1 = 1, X_2 = 1$ の場合

$$Y_i = \beta_1 * 1 + \beta_2 * 1 + \alpha + \epsilon_i$$

$$Y_i = \beta_1 + \beta_2 + \alpha + \epsilon_i$$

- この場合,  $X_1$  と  $X_2$  が影響する場合の値を示していることになります.
  - ex. 既婚で, 子どもがいる人の幸福度がわかる.

▶  $X_1 = 0, X_2 = 0$ の場合

$$Y_i = \beta_1 * 0 + \beta_2 * 0 + \alpha + \epsilon_i$$

$$Y_i = \alpha + \epsilon_i$$

- この場合, 全ての要因が影響しない場合 (何らかの基準となる点) の値を示していることになります.
- ex.未婚で, 子どもがいない人の幸福度がわかる.

## 仮説を立てる

さて、それでは仮説を立ててみましょう。今回分析するテーマは「主観的幸福度(SUB\_HAP)が子の有無(CHI)と結婚(MAR)によって異なる」かどうかを分析します。二要因分散分析（交互作用なし）の場合は以下のような仮説を立てます。

- ▶ 対立仮説1：主観的幸福度の平均値は結婚によって異なる
- ▶ 対立仮説2：主観的幸福度の平均値は子どもの有無によって異なる
- ▶ 帰無仮説1：主観的幸福度の平均値は結婚によって異なるとはいえない
- ▶ 帰無仮説2：主観的幸福度の平均値は子どもの有無によって異なるとはいえない

これらの仮説のもとに分析を行ないます。

## 平均値をプロットする

さて、最初のお約束です。平均値をプロットしましょう。まずは各自でやってみましょう。

さて、例によってggplotguiを使いましょう。

以下のコードはConsole（コンソール）に直接打ち込みます。

```
library(ggplotgui)  
ggplot_shiny(exdataset)
```

そうすると新しいウィンドウが開きます。

以下の通りの作業をしましょう.

- ▶ ggplotタブへ
- ▶ "Type of graph:"は"Dot + Error", Y-variableは"SUB\_HAP", X-variableは"MAR"を設定
- ▶ "Group(or colour)"をCHIに変更
- ▶ "Confidence Interval:"を95%にする.
- ▶ R-codeタブへ行って, 以下のコードのうち, 真ん中のみを以下にする. -また, コード内のdfをexdatasetに変える.

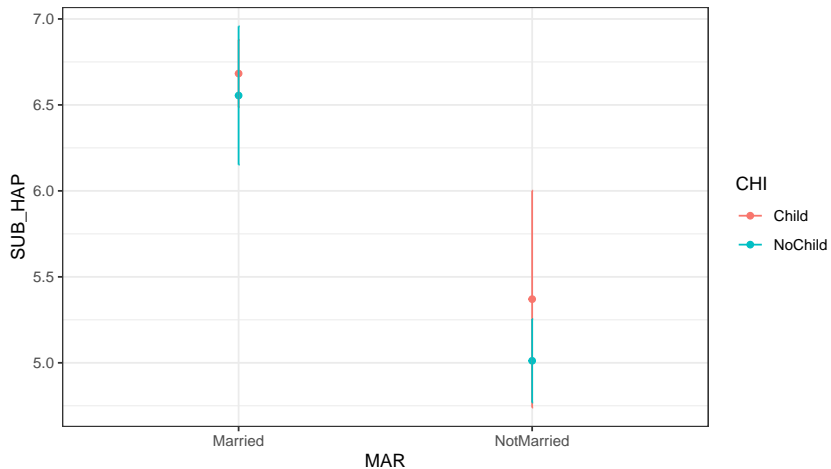
*# You need the following package(s):*

```
library("ggplot2")
```

*# The code below will generate the graph:*

```
graph <- ggplot(exdataset, aes(x = MAR, y = SUB_HAP, colour =  
  geom_point(stat = 'summary', fun.y = 'mean') +  
  geom_errorbar(stat = 'summary', fun.data = 'mean_se',  
               width=0, fun.args = list(mult = 1.96)) +  
  theme_bw()
```

graph



このグラフを見る限り，未婚者に比べて既婚者の方が主観的幸福度が高そうですが，子の有無の影響はありそうな気がしますし，なさそうな気がしますし何とも言えません．したがって，この点についても統計的に差があるのかどうかを明らかにしましょう．

```
marchihap_model_noint <- lm(SUB_HAP ~ MAR+CHI, data = exdata)
# marchihap_model_noint
```



```

#
summary(marchihap_model_noint)

##
## Call:
## lm(formula = SUB_HAP ~ MAR + CHI, data = exdataset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.6991 -1.6991  0.3009  1.3009  4.9667
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.6991     0.1015  66.015 < 2e-16 ***
## MARNotMarried   -1.4642     0.1862  -7.863 1.01e-14 ***
## CHINoChild      -0.2016     0.1832  -1.100   0.271
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1

```

## 結果の書き方

交互作用のない分散分析により、主観的幸福度は結婚しているか否か、および子どもがいるか否かによって異なるかを分析した。その結果、結婚の影響は $F(1, 960)=120.71(p<.001)$ であり、結婚は主観的幸福度に対して有意に影響を与えることが明らかとなった。一方、子の有無の影響は $F(1, 960)=11.21(p>.05)$ であり、有意な影響は認められなかった。

## 結果の解釈

この結果は以下のように解釈することができます.

$$(SH) = 1.464(Married\_dum) + 0.202(Child\_dum) + 5.033$$

ただし, 以下のように変数を割り振っています.

結婚: 未婚→0, 既婚→1

子ども: 子なし→0, 子あり→1

したがって, 結婚と子の有無の影響は以下のように表すことができます.

## 結果の解釈

- ▶ 「未婚者かつ子なし」

$$(SH) = 1.464 \times 0 + 0.202 \times 0 + 5.033$$

$$(SH) = 5.033$$

- ▶ 「未婚者かつ子あり」

$$(SH) = 1.464 \times 0 + 0.202 \times 1 + 5.033$$

$$(SH) = 0.202 + 5.033 = 5.235$$

## ▶ 「既婚者かつ子なし」

$$(SH) = 1.464 \times 1 + 0.202 \times 0 + 5.033$$

$$(SH) = 1.464 + 5.033 = 6.497$$

## ▶ 「既婚者かつ子あり」

$$(SH) = 1.464 \times 1 + 0.202 \times 1 + 5.033$$

$$(SH) = 1.464 + 0.202 + 5.033 = 6.699$$

ここから、未婚者に比べて既婚者の主観的幸福度が高いことはわかりますが、子の有無は主観的幸福度に対して影響をどうも与えなそうです。

## モデル選択：

モデル選択とは、複数の統計モデルを比較する時に用いる手法です。ここでは、モデル選択の手法として、分散分析によるモデル選択とAICに基づくモデル選択を紹介します。

## 尤度比検定によるモデル選択：

分散分析に基づいた近似計算とは，2つのモデル式をもとにして，分散分析を用いることでモデル選択をすることができます．

ここでは，主観的幸福度を応答変数として，説明変数として未婚と子どもの有無を設定したモデルについて，交互作用ありとなしの2つを比較します．

```
anova(marchihap_model, marchihap_model_noint)
```

```
## Analysis of Variance Table
```

```
##
```

```
## Model 1: SUB_HAP ~ MAR * CHI
```

```
## Model 2: SUB_HAP ~ MAR + CHI
```

```
##      Res.Df      RSS Df Sum of Sq      F Pr(>F)
```

```
## 1      959 4695.7
```

```
## 2      960 4697.4 -1    -1.6883 0.3448 0.5572
```

この結果はモデル1である交互作用ありモデルと、モデル2である交互作用なしモデルを比較すると、どちらのモデルも差があるとはいえない確率が55%もあるということを示しております。

この場合は、より単純なモデルとして交互作用のないモデルを選択します。

ちなみに、「分散分析」は実は作成したモデルと説明変数の入っていない「ヌルモデル」を比較しているものと同値になります。



## AIC :

AICとはAkaike's Information Criterion (赤池情報量規準)と呼ばれるものであり, モデル評価の規準の一つです.

$$AIC = -2\log(\quad) + k \times (\quad)$$

として算出され, この値が最小になるモデルを採択します. 特に, 2つのモデルを選択する時には2つのモデルについてAICの差分が2以上あるとそのモデルを選択することができます.

まずはAICを算出してみましょう.

```
AIC(marchihap_model,marchihap_model_noint)
```

```
##                df        AIC
## marchihap_model      5 4268.608
## marchihap_model_noint 4 4266.954
```

ここでは交互作用なしのモデルの方が小さい値を示しています。2つのモデルのAICの差分が2はギリギリありませんが、この場合は交互作用なしのモデルを採択してもよいかと思えます。s