

## 研究ノート

## 意味論を加味した説明変数の選択

——修正済み決定係数  $Q^2$  の改善手法 ——川 瀬 友 太<sup>1)</sup>  
平 井 孝 治

## 目 次

- § 1 重回帰分析に向けた解析枠の確定  
§ 2 t値など各パラメーターとの関係  
§ 3 修正済み決定係数  $Q^2$  の改善手法  
終わりに

## § 1 重回帰分析に向けた解析枠の確定

## 1. 解析枠を確定するための留意事項

昨今、ExcelやSPSS等の解析ソフトの充実により、様々な場面で主成分分析や重回帰分析等の多変量解析が用いられている。その用途としては、大学における授業評価や病院経営実態調査など、多種多様な分野が考えられる。これら調査変数が多い場合には、多変量解析を行うことにより、単純集計だけではわからない、調査項目間の関係や知見を導き出すことができる。

このように多変量解析を行うことにより、豊かな知見を得ることができるが、どのような調査においてもその設計段階において、解析手法を検討しておくべきである。そこで重要になってくるのが、「理論」と「経験」と「思想」の3つである。そこで多変量解析を行う際には、その前提となってくる「理論」を取得しておく必要がある。いかなる手法であれそれを学ぶ際には、その手法の土台となっている「理論」を弁えておくことが作法というものであろう。

次に重回帰分析の精度を高めるためには、実際に「経験」し、解析を行っていく過程において、手法を身に付けていくことが肝要である。因みに、本学経営学部の平井研究室では、調査を経験していく毎に新手法を開発しており、今回のような意味論を加味した説明変数の選択に関する論文を書くに至った。

作法の3つ目として、アンケート調査などで調査票を作成する際に、予め調査対象者の位置付けを行い、調査の結果どのような知見を得たいのか、「出力設計」の重要性を述べておきたい。特に重回帰分析を行う際には、単に決定係数の高いモデルを作ればよいのではなく、説明変数を選択する際に調査主体の「思想」が反映されていなければ、調査目的の達成を期し難

---

1) 立命館大学大学院 経営学研究科 博士課程前期課程 1 回生

いと考えている。

本章では具体的に、重回帰分析を行う際の前提となる、「解析枠の整備」を中心に述べる。

## 2. 優先3択など個数を限定した複数選択項目の回避

調査票の作成時、回答者の意見を聞く場合などに、複数選択の設問を設けることがある。その際に、多変量解析に馴染むようにそれぞれの変数における数値・数量化を検討しておかなければならない。

複数選択の設問に対する数値化は、大きく分けて以下の3つのケースが挙げられる。①個数で持って、当該サンプルの回答者の設問に対する数値・数量とする「個数処理」、②選択肢を個数制限のない、それぞれ独立した変数(調査項目)とみなし、選ばれた変数のみに「1」を入力する「0, 1 処理」。または持ち点(例えば6点)を振り分けていく「持ち点処理」、③選択肢を選ぶ際、選択個数に制限を設け、その個数に従って、「0, 1 処理」や「持ち点処理」を行う場合である。

ここでは選択肢を制限する設問である「優先3択」を例に、重回帰分析をおこなう際の問題点を検討しておきたい。

(Q) 仕事を選択する上で、何が重要であると考えますか。次の中から上位3つまで選択して下さい。

- |                   |            |             |
|-------------------|------------|-------------|
| 1. 能力を活かせること      | 2. 給与・賃金   | 3. 勤務地      |
| 4. 興味があること        | 5. 社会貢献    | 6. 休日・勤務時間  |
| 7. 組織に縛られないこと     | 8. 安定した生活  | 9. 社会的ステイタス |
| 10. 新たな課題に挑戦できること | 11. その他( ) |             |

上記の優先3択の設問は、選択できる変数が3つに限定されているため、選ばれた変数とさもない変数が排他的になり、その変数間で負の相関がおこる。これは重回帰分析による線形モデル構築の際に、障害になる。

このような場合1～10までの各選択肢を独立した変数とし、第1位に選ばれた選択肢には3点、2位には2点、3位には1点を付与して処理するのが通常である。あるいは単純3択でも個数に制限を付ければ負の相関が生ずるので、該当するものを全て選んでもらう全択方式にし、「0.1 処理」をお奨めしたい。この操作によって、説明変数間に生ずる負の相関を和らげ、いい重回帰をもたらす効果が生ずる。

## 3. (0, 1) 変数で、一方が5%未満の場合

②において紹介した「0, 1 処理」においても、重回帰分析を行う際の注意を要する。当該変数の単純集計において、一方が5%未満の場合には、重回帰モデルを構築する際に、想定外の挙動をすることがある。そこで、このような調査変数は、当初から説明変数に組み込まないことをお奨めしたい。

## 4. 2つの解析枠を橋渡しする接続変数

ここでは、異なった母集団で、共通した調査変数を設ける場合を論ずる。一方の解析枠で構築した線形モデルを使って、他方の解析枠に属するサンプルを説明することができる。その例として、ある大学の卒業生と在学生の二つの異なった母集団における調査を紹介したい。卒業生のキャリアに関連する目的変数を、在学中に關係する説明変数だけで重回帰分析したとする。

$$u_i : \text{キャリア変数 (実際値)} \quad x_i, y_i, z_i : \text{在学中変数} \quad a, b, c : \text{回帰係数}$$

$$\text{回帰値} : \hat{u}_i = ax_i + by_i + cz_i + k \quad k \text{ は定数項} \quad (i = 1, 2, \dots, \text{卒業生ナンバー})$$

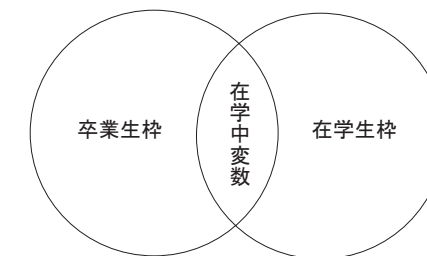
卒業生枠における重回帰で、各変数に対する回帰係数が求められる。求められた回帰係数を在学生の枠に外挿することにより、在学中変数の如何により、在学生の卒業後のキャリアマップをある程度描くことができる。

$$\text{回帰値} : \hat{u}_j = ax_j + by_j + cz_j + k \quad (j = 1, 2, \dots, \text{在学学生ナンバー})$$

これにより、在学学生が現在どのような学び方をしているかで、卒業後どのようなキャリアを築いていけるかを、ある程度予測することができる。

そのためには、予め在学中にどのような学び方をしていたか、卒業生と在学学生に共通した在学中の変数群を組み込んでおくが必要になる。以下に、接続変数のイメージ図を載せておく。

接続変数のイメージ図



## 5. 説明変数群の論理

重回帰分析において、数ある変数の中から、モデル構築に貢献する説明変数を選び出すことが必要となる。当該回帰式(モデル)が目的変数をどの程度説明できているかを示す指標とし

て、決定係数が挙げられる。全変数で重回帰した時に、決定係数  $R^2$  が 0.490 以上ならば、自由度修正済み決定係数  $Q^2$  が 0.4096 以上のモデルが得られる可能性がある<sup>2)</sup>。後者の基準値をこのように設定する理由については、後に説明する。

説明変数を選択する際に、調査主体に「理論」や「経験」が備わっていれば、上記のような水準をある程度担保することができる。しかし、いかに信頼のおけるモデルを構築することができて、そこに設計者の「思想」が反映されていなければ、豊かな知見を得ることは期し難い。

実際の変数選択の作法については § 3 に詳しいので、そちらに譲る。

## 6. サンプル数 $n >$ 調査変数 $p + 1$

重回帰分析を行う際には、サンプル数  $n$  を変数の数  $p$  より多く確保しておきたい。その理由としては、重回帰分析を行う前提となる回帰平面を張ることができないためである。さもない場合は安定した線形モデルを確定できないため、解析を行うことが困難となる。

また 2 つ目の理由としては、行列の演算により導ける正方行列の次数は、さもない場合にはサンプル数に依存することが挙げられる。そのため、各変数に対する共分散や相関行列を導き出せない。変数をサンプル数と同数にすることにより対処できるが、得られる回帰式は一意的に定まってしまう、自由度不足に陥ってしまう。

3 つ目の理由としては、重回帰分析のモデルがどの程度目的変数を回帰できているかを示す指標である、自由度修正済み決定係数  $Q^2$  が問題となる。

$$n : \text{サンプル数} \quad p : \text{変数の数} \quad R^2 : \text{決定係数}$$

$$Q^2 = 1 - \frac{n-1}{(n-1)-p} (1 - R^2)$$

この定義式より、 $Q^2$  はサンプル数  $n$  が説明変数の数  $p$  より小さければ、定義できない。これでは当該モデルが  $Q^2$  を求める前提条件を満たしていないために、どの程度信頼できるかを確かめられない。それゆえ、以下の議論では  $n \geq p + 2$  を仮定する。

## § 2 $t$ 値など各パラメーターとの関係

### 1. モデル内における説明変数の有意性

重回帰分析において、説明変数の当該回帰式に対する寄与の度合いを統計的に計る、 $t$  値、 $p$  値、 $F$  値なる検定量が存在する。

### 1) 母集団 $\Omega$ に想定される「Virtual に真な回帰式」

本節では、「偏差、残差、誤差」の三差がそれぞれ固有の意味を有する。偏差とは  $x_i - \bar{x}$  のことで、残差とは実際値と回帰値との差  $e = u_i - \hat{u}_i$  のことである。

よって以下では、各データを  $x_i : x_i - \bar{x}$  で置換済みとして論ずる。母集団  $\Omega$  からサンプリングして得られたサンプル数  $n$  の標本空間の下で、説明変数が 3 つだとしても一般性を損なわない。

回帰式  $\hat{u}_i = ax_i + by_i + cz_i \dots\dots (i)$  が得られたとすると

実際値は  $u_i = (ax_i + by_i + cz_i) + e_i \dots\dots (i)'$  となる。

しかし、(i) は、母集団  $\Omega$  に属する全数  $N$  で回帰して得られる Virtual に真な回帰式

$$\hat{u}_i = ax_i + \beta y_i + \gamma z_i \dots\dots (*) \quad \text{そのものではない。}$$

しかし、標本空間から得られた回帰係数  $a$  は期待値  $\alpha$  の周辺であるバラつき  $\sigma_a^2$  でもって正規分布していることが知られている。即ち、あくまでも (i) 式は (\*) 式の模造品に過ぎず、サンプルの取り方によって回帰係数  $a, b, c$  は異なってくるが、しかし、 $a$  は正規分布  $N(\alpha, \sigma_a^2)$  に従うことが知られている。このときの、 $a$  と  $\alpha$  の差を「誤差」というが、 $b, c$  についても同様である。

然るに  $a$  の分散  $\sigma_a^2$  は Virtual に真な回帰式を用いて

$$\text{実際値 } u_i = (ax_i + \beta y_i + \gamma z_i) + \epsilon_i \dots\dots (*)'$$

で規定される残差  $\epsilon_i$  の分散  $\sigma^2$  に依存することが知られている。

しかし、この  $\sigma^2$  は一般に不明なので、その代わりにモデル (i)' の残差変動を自由度  $n - (p + 1)$  で修正した残差分散  $\sigma_e^2$  を用いて代用・代理する。

$$\text{その結果、} \sigma_a^2 = \text{収差分散} \times \frac{\text{調整係数 } k_x}{x \text{ の全変動}}$$

その平方根  $\sigma_a$  を回帰係数  $a$  の「標準誤差」と称し、

「帰無仮説  $H_0$  : 回帰係数」の検定に用いる。

### 2) 標準誤差

サンプルサイズ  $n$  の下で、目的変数  $u$  も 3 つの説明変数  $x, y, z$  とも偏差済みで計算された重回帰式で

$$\text{回帰値 } \hat{u}_i = ax_i + by_i + cz_i \dots\dots \textcircled{1} \quad \text{とする}$$

これは元来、当該母集団に存在する (Virtual に真な) 回帰式

$$\hat{u}_i = ax_i + \beta y_i + \gamma z_i \dots\dots \textcircled{2} \quad \text{に対する推定結果である。}$$

例えば、推定係数は回帰母係数  $\alpha$  の偏りのない推定値 (不偏推定値) で、

その期待値は  $E(a) = \alpha$  である。

2) 本論文では、自由度修正済み決定係数を  $Q^2$  と表す

更に例えば、説明変数が  $x$  と  $y$  だけに限った際に、 $x$  と  $y$  の相関係数を  $r$  とすると、分散は

$$V(a) = \sigma^2 \frac{(1 - r^2)}{\sum (x_i - \bar{x})^2} \dots\dots\dots ③ \quad \text{となる。}$$

ここに、 $\sigma^2$  は 実際値 :  $u_i = \alpha x_i + \beta y_i + \epsilon_i \dots ④$  の誤差項  $\epsilon_i$  の分散である。

$$\text{即ち、} \sum \epsilon_i = 0 \text{ で } \sigma^2 = \frac{1}{n} \sum \epsilon_i^2 \dots\dots\dots ⑤$$

なお、 $\sqrt{V(a)}$  を推定値  $a$  の「標準誤差」という。

ただし、 $\sigma^2$  は実際には不明なので、 $e_i = u_i - \hat{u}_i$  とすると

$$V(e) = \frac{\sum (e_i - \bar{e})^2}{n - (p + 1)} = \frac{\sum e_i^2}{n - (p + 1)} \dots\dots\dots ⑥ \quad \text{でもって代理 (surrogate) する。}$$

この重回帰の場合には説明変数の数は  $p = 2$  である。このとき、サンプルサイズがある程度大きい ( $n > 121$ ) ならば<sup>3)</sup>、回帰母係数  $\alpha$  の 95% 信頼区間は

$$\alpha \sim a \pm 1.96 \sqrt{V(a)} \dots\dots\dots ⑦ \quad \text{となる}$$

### 3) 回帰母係数 $\alpha$ の検定

回帰母係数  $\alpha$  は当該説明変数  $x$  の目的変数  $u$  への寄与を表すので、統計的に検定しておく必要がある。即ち、計算された  $\alpha$  の値が「ゼロではなく本当に目的変数に寄与しているか否か」その有意性を検定する。

以下、この「ゼロ仮説の検定」をサンプルサイズ  $n$  がさほど大きくないとして ( $t$  検定) で紹介・議論する。

### 4) 重回帰式 $t$ の値 $p$ と値と $F$ 値

帰無仮説  $H_0$  : 回帰母係数  $\alpha = 0$  を設けて (即ち、 $N\left(0, \frac{\sigma^2}{V(a)}\right)$  の下で)

検定量 :  $t = \frac{\alpha}{\sqrt{V(a)}} \dots\dots\dots ⑧$  を計算する。

この値を当該説明変数の  $t$  値というが、この  $t$  値が大きいほど回帰式に対する寄与が大きいということになる。更に、 $t$  分布関数の逆関数であり、 $\alpha = a$  だとしても、それで誤る確率が  $p$  値である。

では、 $t$  値がどの程度高ければ、採用すべきなのかは、当該データのサンプル数や調査変数に

依拠するが、自由度  $f$  が  $f > 121$  の際、 $t$  値が 2.0 以上であれば当該を採用すべき説明変数と見做してもよい。即ち、自由度  $f > 121$  では、正規分布に近似しているので、信頼度 95% を考えると、両側の検定量の  $1.96\sigma$  とほぼ同値である。なお、 $F$  値は  $t$  値の平方をとった検定量である。

### 5) 多重共線性

重回帰分析における多重共線性とは通称「マルチコ」と呼ばれるものであり、当該重回帰モデル内では不安定な説明変数であることを示している。また、目的変数と説明変数との 2 変数間での単相関係数と、モデル内での偏相関係数の符号が正負反転することがある。そのため、重回帰分析において、モデルの完成度を上げるために、マルチコを排する選択作法を身に付けておくことを誰しも首肯せざるを得ない。更に、その作法を取得するためにも、重回帰分析に関する理論的造詣を深めておかなければならない。

ここに

$$\hat{u}_i = \alpha x_i + \beta y_i + \epsilon z_i$$

とした時、もし

$$z_i = \alpha x_i + \beta y_i$$

という線形関係があったとすると、例えば、0.4 ゆずって  $z_i$  に

$$0.4(\alpha x_i + \beta y_i) + 0.6\epsilon z_i \text{ を代入すると、}$$

$$\hat{u}_i = (a + 0.4\alpha) x_i + (b + 0.4\beta) y_i + 0.6\epsilon z_i$$

となり、回帰係数の符号が変わる可能性がある。

そこで、その線形関係を明らかにするために、単相関と偏相関に着目する。

### 6) 偏相関

偏相関とは、「重回帰モデルにおいて、特定の説明変数が、その説明変数抜きで説明できなかった部分を、どの程度説明・回帰できているのか」を判明させる指標である。偏相関は複雑な指標であり、概して、当該説明変数が他の説明変数との関係で、どの程度モデルに貢献しているかを示すものと考えてよい。目的変数と説明変数の単相関が 0 の場合でも、偏相関 (の絶対値) が 1 近くになることもある。

### 2. $Q^2 \geq 0.4096$ に設定する理由

自由度修正済み決定係数  $Q^2$  がどの程度あれば、モデルの名に相応しいのかは、解析の目的や得たい知見にもよる。しかし、一般に明確な数学的基準が無いために、解析主体が信頼に値すると主張する重回帰モデルの精度が疑わしい場合も多々ある。

そのため、自由度修正済み決定係数  $Q^2$  に明確なボーダーラインを設け、モデルの精度を担

3) ある程度大きいということ、本学経営学部平井研究室では  $n > 121$  としている。

保し、知見を豊かにする価値は大いにある。また、当該モデルの信頼性を確保するために、どの程度自由度修正済み決定係数  $Q^2$  の値を保持するのか、解析者が論拠を明示する必要がある。筆者らが考える  $Q^2 \geq 0.4096$  の論拠は、以下の 3 点である。

### 1) 0.64 の平方

ピタゴラスの定理で斜辺を 1 とすると、代表的な直角三角形のそれぞれの辺の比は、1 : 0.8 : 0.6 である。そこで、しばしば、0.8 の平方の 0.64 を強い相関の下限値とし、0.6 の平方の 0.36 を弱い相関の下限値とすることがある。そこで、この 0.64 の平方を求め、0.4096 を新たに境界線 (Boundary) として設定する。これが、 $Q^2 \geq 0.4096$  に設定する 1 つ目の理由である。

### 2) 1.64 $\sigma$

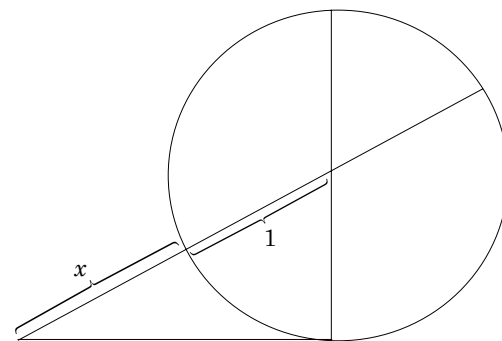
統計学では、正規分布を用いて推定や検定をすることが多々ある。その際、頻繁に用いられる検定量は 1.64 $\sigma$  と 1.96 $\sigma$  であり、それぞれ、両側検定だと信頼係数が 90%、95% となる値である。検定量の 1.64 を  $2^2$  で割ると、0.410 であることから、 $Q^2 \geq 0.4096$  の評価基準を設定する 2 つ目の理由とする。ちなみに他方の 1.96 を  $2^2$  で割った値は、0.49 で我々が修正前の決定係数  $R^2$  の基準値にしているものである。

### 3) 黄金比 (Golden Ratio)

比率の中で、人が最も美しいと感じられるといわれている、黄金比 (Golden Ratio) がある。自然界に頻繁に現れる比率であり、歴史的建造物にも用いられている。黄金比は、半径 1 の単位円を仮定したときに右図のように、円の中心を通る線分を引き、それぞれの長さを  $x$  と 1 と置く。このとき、 $x : 1 = x + 1 : x$

となり、比は  $\frac{\sqrt{5}-1}{2} : 1$  もしくは  $\frac{\sqrt{5}+1}{2} : 1$  となる。これを黄金比としている。

そこで、弱い相関の下限値として用いられている、0.6 の平方の 0.36 と、0.7 の平方の 0.49 ( $R^2$  の下限値) をそれぞれ始端と末端にした、数直線を考える。この 2 点の間を黄金比  $\frac{\sqrt{5}-1}{2} : 1$  で分割すると、その分割点 (Golden Section) は 0.4096 となる。これが、1) 2) で設定した値に限りなく近づき、 $Q^2 \geq 0.4096$  を設



けた 3 つ目の理由となる。

### 3. $Q^2 \geq 0.4096$ に設定する効用

$Q^2 \geq 0.4096$  に設定する効用には以下の 3 つが考えられる。

#### 1) 再現性の担保

誰もがその手法・方法を用いれば、同様の結果が得られるという再現性は科学する際の重要な要件である。この  $Q^2 \geq 0.4096$  を目指して、説明変数を選択していくことで、解析者が誰であれ、同様、もしくは近似した精度の結果が得られることが判る。即ち、変数選択の作法を身に付けていれば、信頼に値するモデルを構築できる可能性があるといえる。その意味で厳しくはあるが、この  $Q^2 \geq 0.4096$  を設定し、それを満たすような解析の過程を踏むことで、モデル構築の再現性が担保されると考えられる。

#### 2) モデルの精度向上

説明変数を  $Q^2 \geq 0.4096$  を目指して、選択していくことで重回帰モデルの値が格段に向上する。例えば、全変数で重回帰したモデルでは、自由度修正済み決定係数  $Q^2$  が限りなく 0 に近い値であるにも関わらず、変数を削減すれば  $Q^2$  が復活し、0.4 前後まで向上することがある。これは、説明変数の選択に「思想」ないし「意味論」を加味したところが大きいと思われる。SPSS を用いた重回帰分析では、解析ソフトのアルゴリズム (例えば、ステップワイズ法) によって導き出されたモデルこそが、優れたモデルと判断しがちである。しかし、そこには意味論が存在せず、解析ソフトのパズルゲームの結果といわれても反論し難いものがあり、ましてや分散比の  $P$  値が 1% 以下だからといって、その程度では決していいモデルとは限らない。解析者が自ら「思想」を持って、挑まなければ、そこから優れた知見は引き出せず、まして調査目的の達成はとうてい困難であろう。

#### 3) 説明変数の数の最小化

$Q^2 \geq 0.4096$  という厳しい基準を満足するためには、必然的に少ない説明変数で、簡明にモデルを構築することが求められる。即ち、回帰式の説明変数を減らし、本質的にモデル構築に資する説明変数を選びきることが必要となる。そのため、多くの説明変数を用いて説明するよりも、因果関係が明らかになる効用がある。「真実は簡明を尊ぶ」と考える作法こそが、モデル構築に資するのである。

#### 4. $Q^2$ が明らかに高い場合の選択

$Q^2$  が極端に低い際、 $Q^2 \geq 0.4096$  のようにある閾値を設けてモデルを構築していくべきである。しかし、サンプル数が十分にある場合などは、 $Q^2$  が高い場合が多く、説明変数をいかように組み合わせても、優れたモデルが構築されるであろう。しかし、如何様に組み合わせたところで「思想」ないし「意味論」が反映されていないモデルは、これまでに述べた理由から一考を要する。

即ち、SPSS などの解析ソフトだけでは設計・解析の思想を反映したモデルを構築することは難しく、①どのような知見を得たいのか、②いかに設計思想を反映させるか、③説明変数群の論理整合性をどう担保するのか など念頭に置き、モデルの構築に挑むべきであろう。

### § 3 修正済み決定係数 $Q^2$ の改善手法

複数ないし単数の目的変数を対象に重回帰分析し、当該分野の線形モデル構築し、分析や仮説の検証に資することがしばしばある。分析や仮説の検証を行う際に、目的変数と説明変数の選択やその因果関係が非常に大きな意味を持つことは、言うまでもない。そこで、本節では説明変数の選択に意味論を考慮したモデル構築の手法について紹介する。なお、目的変数の選択については、立命館経営学第 45 巻第 2 号「目的変数の合成に関する課題の考察—病院における人事評価を例として—」山本友太ほか 著 を参照されたい。

#### 1. 意味論を考慮したモデル構築

なぜ、意味論を加味して説明変数の選択を行うのか、それは以下の 3 つの理由による。

##### 1) モデルの独立性

重回帰分析は、説明変数の数に制限がかかっているわけではない。よって、説明変数の数やウェイトの小さい変数をほかの変数と入れ替えるだけで、いくつものモデルができる可能性がある。そこに、意味論を加味して説明変数を選択し、モデルを構築する意義がある。

即ち、決定係数  $R^2$  や自由度修正済み決定係数  $Q^2$  を高めるために、関係のない変数を排除し、解釈するに値しないモデル構築を回避することが出来る。説明変数を安易に増やすと、然るべき因果関係がないにも関わらず、当てはまりの良さの指標である自由度修正済み決定係数  $Q^2$  が上がり、モデルの解釈を誤ることがある。特に少ないサンプル数のデータで、十個ほどの説明変数を採用することにはかなり問題がある。そのため筆者らは、自由度修正済み決定係数に  $Q^2 \geq 0.4096$  という基準値を設けている。だからこそ説明変数の組み合わせには、統計的な重要性と共に意味論を考えながら、変数の数を最小限にすることが求められる。

しかし一方で、サンプル数が非常に多い場合には、どんな説明変数を組み込んでも  $Q^2$  が優に  $0.4096$  を超える場合がある。この場合においても、意味論をもって説明変数を削り込まな

い限り、幾通りものモデルが存在することになる。

よってここで紹介するモデル構築の作法は、 $Q^2 \geq 0.4096$  という基準値を超え、説明変数に意味を持たせた上で削りこみ、変数の数が最小限になるように、削り込む手法である。

#### 2) 説明変数選択のしやすさ

説明変数を選択する際に、ごく一般にステップワイズ変数選択 (step-wise selection) 法を用いる場合が多い。これは、説明変数の候補から、予測や判別に有用な順に独立変数を採用するための方法で、意味論とは無関係であることはいまでもない。例えば、最も有用な説明変数を 1 個採用する。次に、まだ採用されていない説明変数のうちで最も有用な独立変数を 1 個採用する。その際に、最初のほうで採用された説明変数も、後で採用された変数との関係で不要になる場合があるので、新たな説明変数の採用の前に、すでに採用された変数を取り除くかどうかをチェックする。つまり、説明変数を 1 個 1 個どちらがよいのかを判断し、変数を選択していくやり方である。

これは有効な方法かもしれないが 100 変数もあるような調査の場合、かなりの労力が費やされる。極論すれば、100 の説明変数があった場合、nCp 通りも確認し続けることになる。そこで、説明変数を選択する際に、やはり意味論を考慮することを提案したい。つまり、意味論を加味して説明変数を選択すると、ステップワイズ法に全く頼ることなくモデル構築ができる利点がある。

たとえば、06 年度大学コンソーシアム京都「京都教育力調査」で作成した重回帰モデル「キャリア効力」<sup>4)</sup> を見ていただきたい。このモデルは調査変数「Q16-5 本質を見抜く能力」、「Q16-6 問題解決能力」、「Q16-7 人モノ活用」の 3 変数を主成分分析で統合し、仕事やキャリアに関する自信と解釈し「キャリア効力」とした。これを目的変数に重回帰分析を行ったものであるが、意味論を加味して考えると、説明変数候補として、キャリアに関する変数はもちろんのこと、在学中の学び姿勢の良さや、各能力の取得する時期の早さなどが考えられる。もちろん全変数で重回帰分析を行い、説明変数を削っていくのであるが、意味論的に関連する変数に注目して削り込んでいく。すると、削り込む過程で候補に立てた説明変数と目的変数の関係が浮き彫りになるという利点がある。

また、筆者らは基準値を自由度修正済み決定係数で  $Q^2 \geq 0.4096$  としているが、 $Q^2$  が  $R^2$  に収束するとしばしば、 $0.4096$  を削り込むことがある。その際にも意味論を加味した手法が有力となる。すなわち、閾値を越えるモデル構築のため、削った変数の復活を考えるのである。意味論を加味して変数復活を期すると  $Q^2 \geq 0.4096$  となった事例に、筆者らはたびたび遭遇

4) 資料 1「キャリア効力」参照

している。

よって、意味論を加味して説明変数を選択することで、ステップワイズ法に頼ることなく、モデル構築が容易になる。そのみならず、目的変数と説明変数間の関係や変数復活にも有効であることをしばしば体験している。

### 3) モデルの解釈の容易性

前述の 2) 説明変数の選択のしやすさ と関わるのであるが、説明変数選択に意味論を加味しているため、出来上がったモデルの分析・解釈がしやすくなる利点がある。意味論を加味するという事は、説明変数間にそれなりのロジックが存在する。当初から、意味論を持ってすれば、出来上がったモデルについて解釈が容易となる。また意味論を援用すると、説明変数間だけでなく、説明変数と目的変数との因果関係を分析することで、仮説の検証も容易になるという利点がある。

またそれに加え、説明変数間に意味論やロジックを求めることで、モデルの誤った解釈を防ぐという効用がある。つまり、意味論から考えて目的変数とかなり相関の強い変数や因果関係の見えない変数は、説明変数に相応しくなく、除外してモデル構築をする必要がある。たとえば、先述した「キャリア効力」において「自己効力感」を説明変数に入れなかった。キャリア効力と自己効力感の相関係数は 0.3921 で筆者らのいう弱い正の相関があるが、それらの意味が同義反復的な部分もあると解釈したため説明変数から除外した。すなわち、完成したモデルの解釈を容易にすることに加え、意味論を加味して説明変数を選択することで、 $A \rightarrow A$  のようにトートロジーなモデルを排除するのである。

以上 3 点から導き出される利点を確保するには、事前解析を丹念に行うことが求められる。重回帰分析の事前解析として、単純集計はもとより、主成分分析、変数クラスター分析を行う必要がある。事前に解析を行うことで、全体の傾向や変数の特徴をつかむことが重要となる。たとえば、4 択の問題で、1 から 4 まで万遍なくばらつく回答もあれば、3 か 4 に回答が集中する場合もある。また主成分分析が持つ双対性の意味の解釈や、変数クラスター分析における意味の解釈も欠かせない。とくにクラスター分析は正負の相関の強いものがクラスターを形成することから、変数同士の関係が大まかにつかめるメリットがある。大まかな特徴を把握した上で、自らの仮説のもとに意味論を加味し、説明変数の選択を行っていく。そうすることで、ステップワイズ法のような機械的な手間を省き、仮説検証に資するデータが得られるメリットがある。

## 2. 説明変数選択のフローチャート

本節では、意味論を加味した説明変数選択を紹介する。またそれをフローチャートに示すと、

図 1 のようになる。なお、番号はフローチャートに対応させている。

### ①目的変数の選択

単独目的変数と統合目的変数の両者がある。統合目的変数の場合、統合するロジック<sup>5)</sup>や、以下次節で紹介する「定義版」に関わるので、そちらを参照されたい。

### ②全変数で重回帰

目的変数に含んだ変数と、自明なモデルにならないように意味論で排除した変数を除いた変数群(全変数)で重回帰分析を行う。

### ③ $R^2 \geq 0.490$

筆者らは  $R^2 \geq 0.490$  以上である場合のみ、説明変数の削り込みを行う。理由は、基準値が  $Q^2 \geq 0.4096$  ということに加え、 $R^2$  と  $Q^2$  の関係は、 $R^2 \geq Q^2$  であり、紹介する手法で変数を削り込んでいくと  $Q^2$  が  $R^2$  に収束するためである。すなわち、全変数で重回帰分析を行って、最低でも  $R^2 \geq Q^2 \geq 0.4096$  とならなければ、変数を削り込む意味がない。

また、全変数で重回帰分析を行った際、 $R^2$  が優に 0.490 を超えているにも関わらず、 $Q^2$  が 0 となることも筆者らはしばしば体験している。キャリア効力やワークアビリティの事例についても、全変数の重回帰分析では  $Q^2$  が 0.30 前後であった。しかし、変数を削りこんでいく中で  $Q^2$  が  $R^2$  に収束し、当該モデルが望ましいレベルに達することができた。

### ④単相関や偏相関で強い・弱いに区分

解析枠にもよるが、筆者らは単相関で  $\pm 0.10$  を境とし、 $\pm 0.10$  未満の弱い相関の変数には印しをつける。偏相関との関わりや  $Q^2$  の改善には相関の強い変数を残すことが有効なため、あらかじめ印しをつけておく。即ち、偏相関が強くない限り単相関の弱い変数を除外することがあるためである。

ところで、なぜ  $\pm 0.10$  かは、明確な根拠はないが、印のつく割合それなりに妥当性があるからである。仮に  $\pm 0.30$  前後とすれば、今回の事例では大半の変数に印しがついてしまう。しかし、あくまで解析枠の特徴によるもので、絶対的基準でないことは留意されたい。

### ⑤マルチコと $F$ 値の小さい変数を削除

マルチコと意味論を加味しながら  $F$  値  $\geq 1.0$  となるように削除する。おおよそこの時点で、

5) 立命館経営学第 45 巻第 2 号「目的変数の合成に関する課題の考察—病院における人事評価を例として—」山本友太ほか著を参照されたい。

F 値が 1.0 未満の変数が半数以上出現するが、なかでも単相関や偏相関の低い変数を中心に削除することをお奨めする。またこの時に、④で行った印しを頼りに、F 値が低く相関の高い変数が多分にありうる。その際には、当該変数を残しておくことが望ましい。

⑥残った変数で重回帰分析

ここまでで残った変数で再度、重回帰分析を行う。

⑦マルチコなし

⑤において、マルチコが発生した変数を削っているにも関わらず、再度マルチコが発生することがある。ここでの判断は迷うが、マルチコの変数を削らなくても、最終的な結果が同じになることがある。マルチコのついた変数が削りこみの過程で、自然と削除されたりマルチコ印が消えることもしばしばある。このように変数削除の判断は難しいが、最後までマルチコ印が残ることもありうるので、早い段階で削除しておくことをお奨めしたい。

⑧  $R^2 \geq 0.410$  の間,  $Q^2 \geq 0.4096$

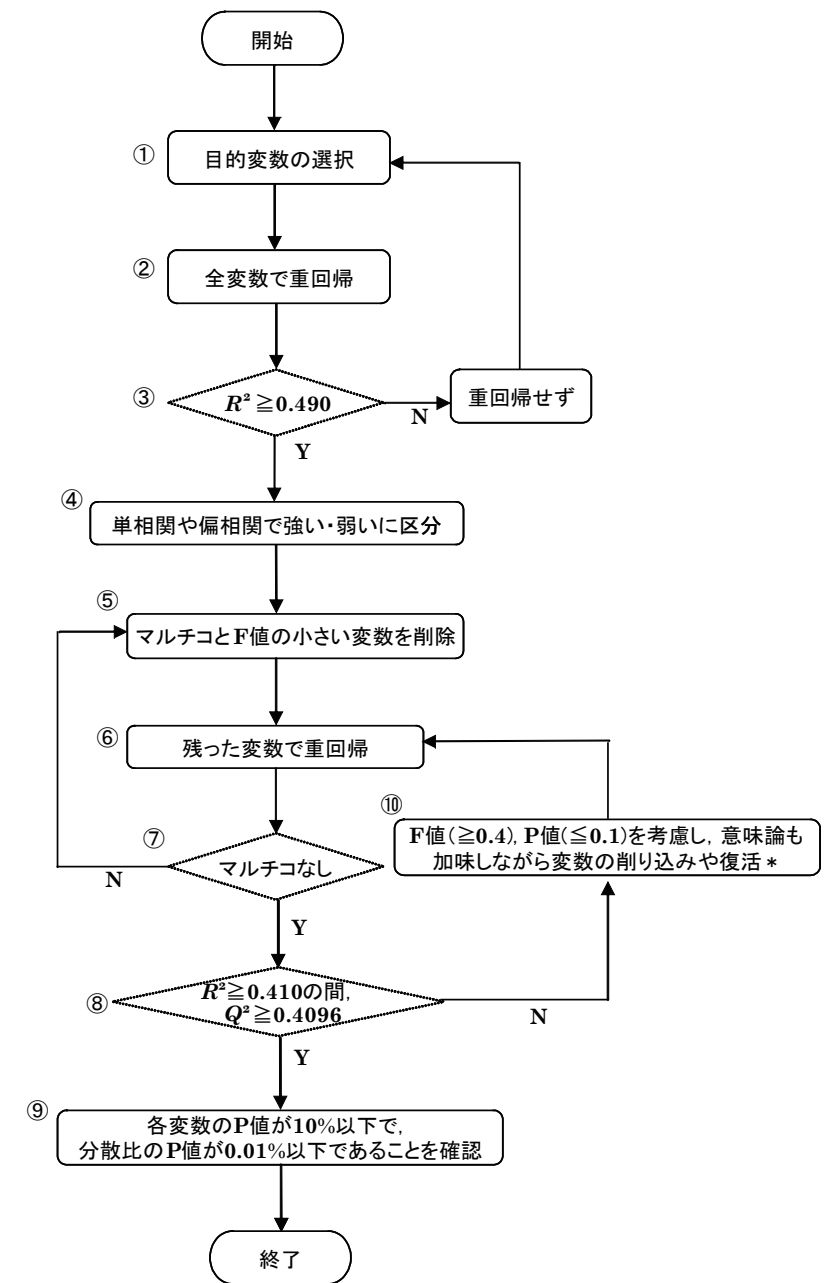
⑦以降は、意味論を加味しつつ ③~⑦もしくは、⑤~⑦を繰り返す、 $Q^2 \geq 0.4096$  になるまで変数を削り込む。その際、あまりにも P 値の高い ( $P \geq 0.1$ ) 変数が残っていたり、マルチコ印つきの変数が残っていたりする場合は、いくら  $Q^2 \geq 0.4096$  であっても、モデルとして採用しない。その際は、単独目的変数なら目的変数の見直しや統合目的変数なら統合する変数に誤りがあると解さざるを得ない。また、そもそもの削り方に問題がある場合もある。いずれにせよ、誤った解釈になるので、モデルとして採用しないことが望ましい。

⑨  $Q^2 \geq 0.4096$

この時点で、 $Q^2 \geq 0.4096$  であること、自明なモデルになっていないこと、意味論的に整合していることなど確認したうえで、当該モデルを採用することをお奨めする。その際、各説明変数の P 値が 10% 以下で、分散比の P 値が 0.1% 以下になっているか否かを確認しておくこと。

⑩意味論を加味し、変数の出し入れ

しかし、 $Q^2 \geq 0.4096$  を下回ることがしばしば起こる。その際に、変数の出し入れを行い、 $Q^2$  を改善させる。変数を削除した際、 $R^2$  がそれほど減らないにもかかわらず、 $Q^2$  が改善することがある。また、過去に削った変数を以下のような方針で復活させることがあるが、⑪~⑬がその際の手法である。



\* 変数の復活とは、主成分分析や変数クラスターなどを加味し、以前に削除していた「変数の復活」を思索

図 1



## ⑪意味論で出し入れ

前節より頻出の「意味論」を加味して、変数の出し入れを行う。その際に、目的変数との意味関係を考えることに加え、説明変数群内でのロジックや因果関係を考慮して、変数の出し入れを行うことをお奨めする。

## ⑫変数クラスターや第一主成分を見て変数の出し入れ

意味論を考える際のヒントになるのが、クラスターや解析枠全体の主成分分析の結果である。特に主成分分析によって得られた第一主成分は、全体の指標になっていることが多く、これをもって目的変数と現時点での説明変数がどのような関係にあるかを考察する。また、その周辺に集まる変数や意味論から考えて、1) 双対な変数ならびに 2) 負の相関を含め強い相関がありそうな変数を見つけ出し、出し入れすることをお奨めする。

## ⑬目的変数と正負の相関の強い変数の出し入れ

⑫でも触れたが、相関係数にも注目する。全変数で行った重回帰分析の結果の考察や、④で印をつけられていない正負とも強い相関をもつ変数を復活させることをお奨めする。

以上をもってして、 $Q^2$  が改善されないならば、単独目的変数なら目的変数の見直しや統合目的変数なら統合する変数に誤りがある場合があるので、①の出发点まで戻って再検討をする必要があるだろう。

## 3. 「定義版」を模索する実験、シミュレーションから実験へ

筆者らは前述のフローチャートに則って幾多の重回帰モデルを作り上げてきた。そして、いくつかの調査変数の主成分による統合によって、当該概念を定義する際、重回帰分析を用いて、模索する実験を行ってきた。筆者らのいう定義版とは、①目的変数群と説明変数群が意図する概念を必要かつ十分に表していること、②説明モデルが  $Q^2 \geq 0.4096$  の条件を満たしているモデルになっていること。つまり、今まで解釈できていなかった概念について本質的な内容を、重回帰分析を行うなかで定義してきた。

⑬で、モデルにならなかったものは、単独目的変数なら目的変数の見直し、統合目的変数なら統合する変数の選択に誤りがあると考えられるので、①まで戻って再検討する必要性を述べた。つまり、モデルにならないということは、統合目的変数が概念を必要かつ十分に表現していないと解釈し、説明変数群と統合目的変数の因果関係や両方の組み合わせを考え、モデルを作ってきた。たとえば、事例の「キャリア効力」や「ワークアビリティ」<sup>6)</sup> がその例であった。

これまでキャリア効力やワークアビリティの概念や因果関係は説かれていなかった。そこで調査変数を用いて、目的変数や説明変数を何度も検討する中で、2つの要件をクリアにするモデルを作り、当該概念を定義した。

つまり、実験の代わりに現象を模擬するシミュレーションとは違い、パソコン上で実験を繰り返す中で定義版の模索を行う手法である。通常、社会科学では研究に実験という手法を用いるには、大変困難をとまったり、長期間を必要とする。しかし、重回帰分析の場合はそれが可能となり、今回のように概念を確定することができた。実験を繰り返し、先にあげた2つの要件を満たすことでモデルができれば、定義として用いることが可能となる。シミュレーションではなく、実験を重ねることで定義を探るこの手法をお奨めしたい。

## 終わりに

本論文は、重回帰分析における説明変数の選択につきソフトに任せることなく、意味論に依拠して線形モデルを構築する方法を主としてスキルに焦点を当て論じてきた。当然のことながら、その背景に数学的な裏づけを取っているが、その理論的な詳細については紙幅の関係上、必要最小限に留めている。

私どもの研究室では、それこそ毎日のごとく誰かが重回帰分析を行っているが、その間に開発した手法も多岐にわたっている。特徴的なことを 2,3 挙げれば、次のようになる。

① 単独目的変数だけではなく、主成分で統合した合成変数を目的変数とすることがある

② 修正済み決定係数  $Q^2$  が 0.4096 以上のモデルを採用すること

③ 意味論を考慮し、無意味な変数や自明な変数をモデルから排除することなどなどである。

これらのほかにも、マルチコの排除などタイトな基準のもとにモデルを構築してきたが、そのおかげで定義版の設定など得られた成果に見るべきことが多々あった。ややもすると、修正済み決定係数の基準切り下げで妥協したいこともしばしばあったが、今では妥協せずに本当に良かったと思っている。一言で言えば、厳しい基準がよい成果をもたらしたということである。フローチャートでも紹介しているが、各説明変数の  $P$  値が 10% を超えないように心がけているが、これは回帰係数の「ゼロ仮説」に耐えるためである。更に分散比の  $P$  値が 0.01% 以下にしているのは、1% 基準ではしばしばおそまつなモデルに出会うからである。

重回帰モデルを構築して、知見検討の段階に入ると予期に反してむしろ目的変数のほうが原因で、説明変数のほうが結果と解したほうが適合する場合もままありうる。そのような場合には、統合目的変数の方をばらして説明変数に置き換え、説明変数群のほうを主成分で統合して目的変数とし、モデルを再構築することもある。このような場合、平井研では先のモデルに対

6) 資料 2 「ワークアビリティ」参照



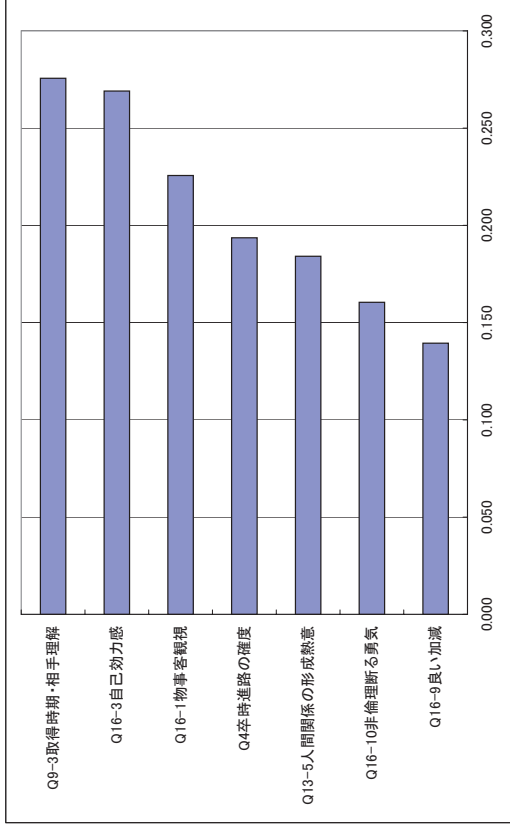


重回帰式 目的変数 ワークアビリティ

Table with 10 columns: 説明変数名, 偏回帰係数, 標準偏回帰係数, F値, P値, 判定, T値, 標準誤差, 偏相関, 重相関, 寄与度. Rows include Q16-9, Q16-10, Q13-5, Q16-1, Q4, Q16-3, Q9-3, and 定数項.

Table with 2 columns: 法定係数, R2, 自由決定係数, R', 重相関係数, R'. Values: 0.469, 0.417, 0.685, 0.646.

Table with 2 columns: 変動, 備置, 自由度, 不備分散, 分散比, P値, 判定. Rows: 全体変動, 帰属による変動, 回帰からの残差変動.

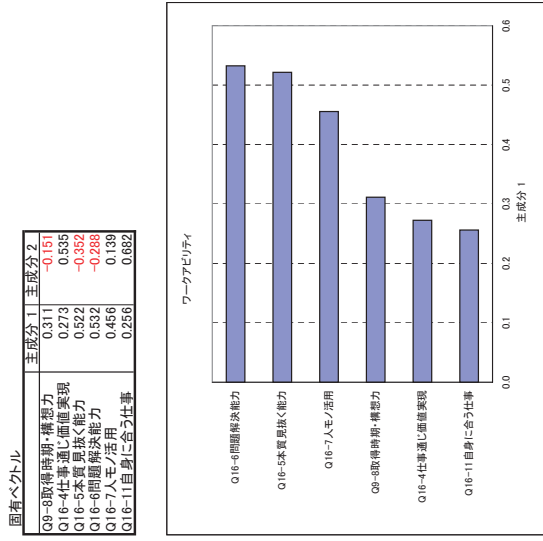


主成分分析 資料 2 合計・平均標準偏差

Table with 4 columns: 合計, 平均, 標準偏差(σ), 標準偏差(σ-1). Rows list variables like Q9-3, Q16-4, Q16-5, Q16-6, Q16-7, Q16-11.

Table with 2 columns: 固有値, 寄与率(%), 累積(%). Rows: 1, 2. Values: 2.54, 1.32, 22.05, 64.22.

Table with 2 columns: 固有ベクトル, 主成分 1, 主成分 2. Rows list variables like Q9-3, Q16-4, Q16-5, Q16-6, Q16-7, Q16-11.

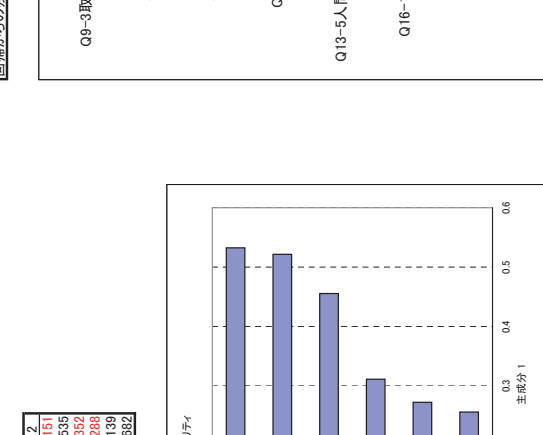


主成分分析 資料 3 合計・平均標準偏差

Table with 4 columns: 合計, 平均, 標準偏差(σ), 標準偏差(σ-1). Rows list variables like Q9-4, Q9-8, Q13-3, Q13-4, Q15-4, Q16-5, Q16-6, Q16-7, Q16-11.

Table with 2 columns: 固有値, 寄与率(%), 累積(%). Rows: 1, 2. Values: 3.36, 1.73, 27.99, 42.43.

Table with 2 columns: 固有ベクトル, 主成分 1, 主成分 2. Rows list variables like Q9-4, Q9-8, Q13-3, Q13-4, Q15-4, Q16-5, Q16-6, Q16-7, Q16-11.

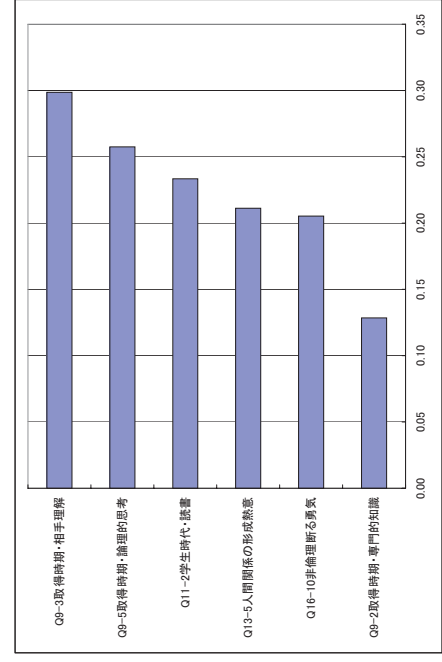


目的変数 生きる力

Table with 10 columns: 説明変数名, 偏回帰係数, 標準偏回帰係数, F値, P値, 判定, T値, 標準誤差, 偏相関, 重相関, 寄与度. Rows include Q9-3, Q13-5, Q11-2, Q9-3, and 定数項.

Table with 2 columns: 決定係数, R2, 自由決定係数, R', 重相関係数, R'. Values: 0.454, 0.410, 0.674, 0.640.

Table with 2 columns: 変動, 備置, 自由度, 不備分散, 分散比, P値, 判定. Rows: 全体変動, 帰属による変動, 回帰からの残差変動.

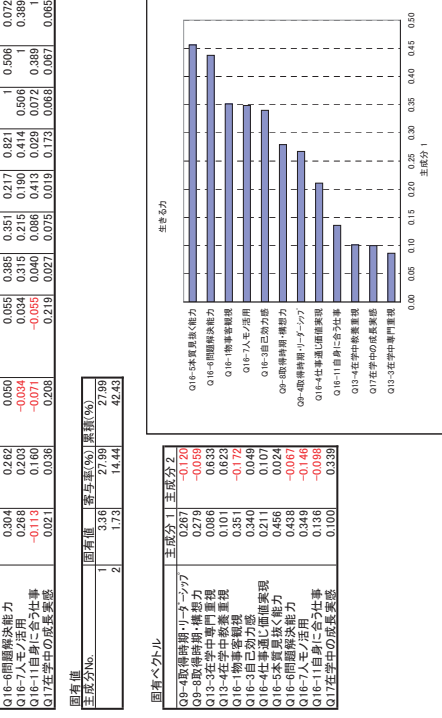


主成分分析 資料 3 合計・平均標準偏差

Table with 4 columns: 合計, 平均, 標準偏差(σ), 標準偏差(σ-1). Rows list variables like Q9-4, Q9-8, Q13-3, Q13-4, Q15-4, Q16-5, Q16-6, Q16-7, Q16-11.

Table with 2 columns: 固有値, 寄与率(%), 累積(%). Rows: 1, 2. Values: 3.36, 1.73, 27.99, 42.43.

Table with 2 columns: 固有ベクトル, 主成分 1, 主成分 2. Rows list variables like Q9-4, Q9-8, Q13-3, Q13-4, Q15-4, Q16-5, Q16-6, Q16-7, Q16-11.



【主成分分析】自己器能

資料 4-1 合計・平均標準偏差

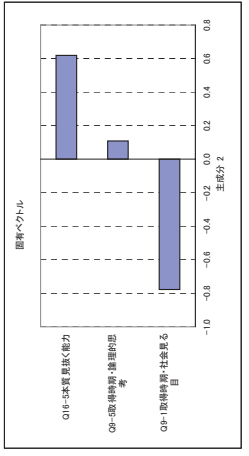
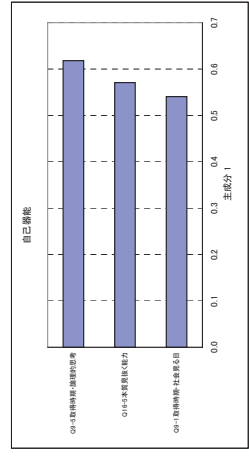
相関行列

因子荷重

主成分 No.

固有ベクトル

主成分 1



【主成分分析】仕事スタイル

資料 4-1 合計・平均標準偏差

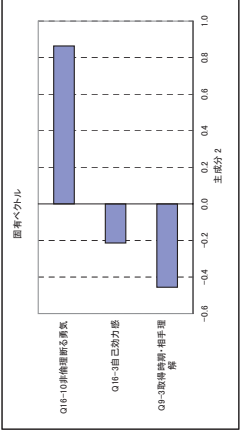
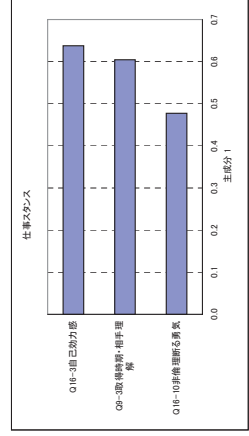
相関行列

因子荷重

主成分 No.

固有ベクトル

主成分 1



【主成分分析】仕事スタイル

資料 4-1 合計・平均標準偏差

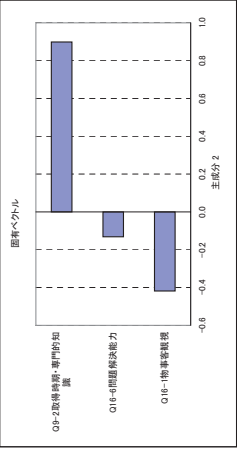
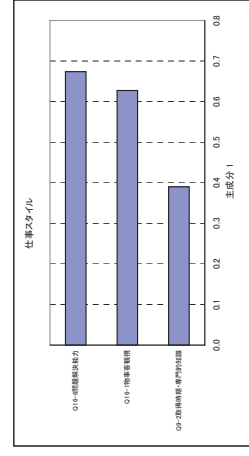
相関行列

因子荷重

主成分 No.

固有ベクトル

主成分 1

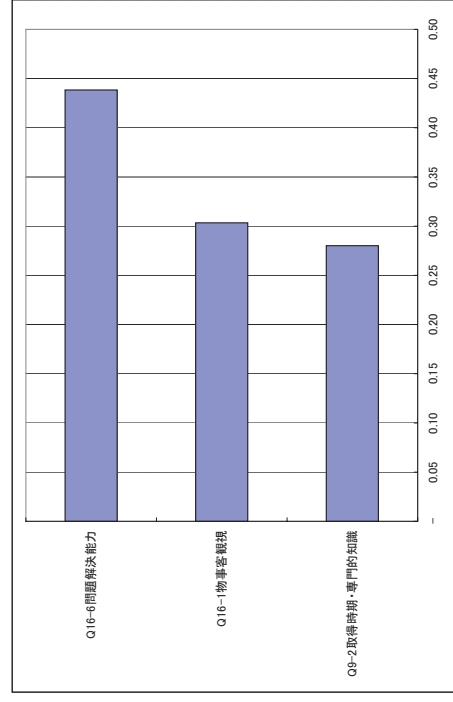


資料 4-2 【重回帰式】自己器能

目的変数

【精度】

【分散分析表】

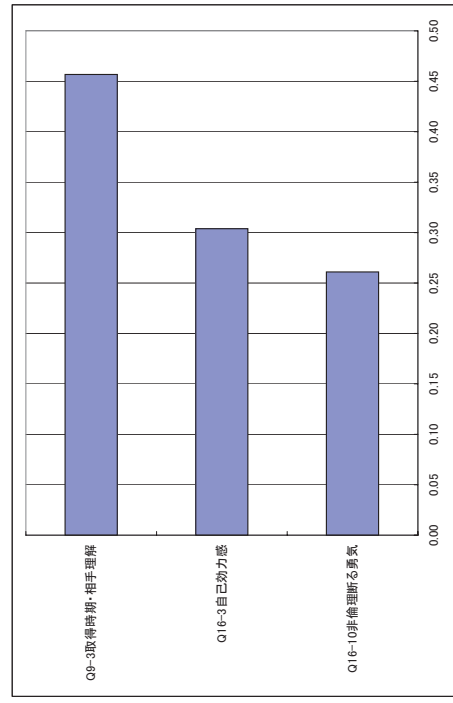


【重回帰式】自己器能

目的変数

【精度】

【分散分析表】



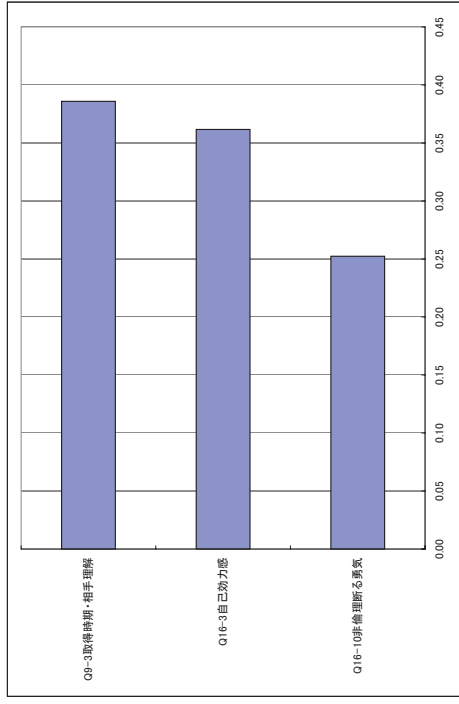
資料 4-3 [重回帰式] 目的変数 仕事スタイル

説明変数名 (自己変数)	偏回帰係数	標準偏回帰係数	F値	P値	判定	T値	標準誤差	偏相関	単相関	符号zz
Q16-10非倫理的な専攻	0.394	0.252	8.482	0.005 [**]	2.912	0.135	0.317	0.334		
Q16-3自己効力感	0.620	0.362	16.848	0.000 [**]	4.105	0.151	0.426	0.472		
Q9-3取得時期・相手理解	0.526	0.366	19.328	3.5E-05 [***]	4.396	0.120	0.450	0.492		
定数項	-3.909									

[精度]  
決定係数 R<sup>2</sup> = 0.441  
自由度数 419  
重相関係数 R' = 0.684  
自由度数 419  
決定係数 R<sup>2</sup> = 0.647  
自由度 419  
決定係数 R<sup>2</sup> = 0.923

[分散分析表]

変動要因	偏差平方和	自由度	F値	P値	判定
全体変動	17.134	79			
回帰による変動	51.13	3	17.044	1.9E-09 [***]	[**]
回帰からの残差変動	64.81	76	0.853		



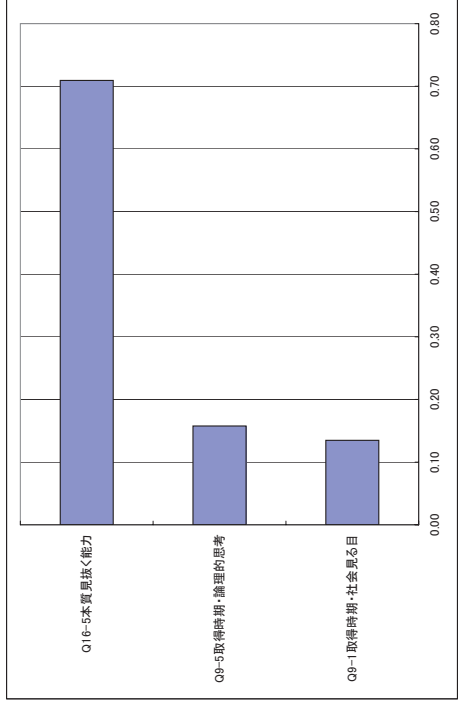
資料 4-4 [重回帰式] 目的変数 仕事スタンス

説明変数名 (自己変数)	偏回帰係数	標準偏回帰係数	F値	P値	判定	T値	標準誤差	偏相関	単相関	符号zz
Q9-1取得時期・社会系志願	0.205	0.135	3.814	0.054 [ ]	1.953	0.105	0.219	0.337		
Q9-5取得時期・論理的思考	0.256	0.158	4.860	0.031 [*]	2.205	0.116	0.245	0.436		
Q16-5本質見抜く能力	1.310	0.709	102.819	9.0E-16 [***]	10.140	0.129	0.798	0.792		
定数項	-4.646									

[精度]  
決定係数 R<sup>2</sup> = 0.678  
自由度数 663  
重相関係数 R' = 0.822  
自由度数 663  
決定係数 R<sup>2</sup> = 0.814  
自由度 663  
決定係数 R<sup>2</sup> = 0.703

[分散分析表]

変動要因	偏差平方和	自由度	F値	P値	判定
全体変動	17.134	79			
回帰による変動	26.118	3	52.816	1.5E-18 [***]	[**]
回帰からの残差変動	37.58	76	0.495		



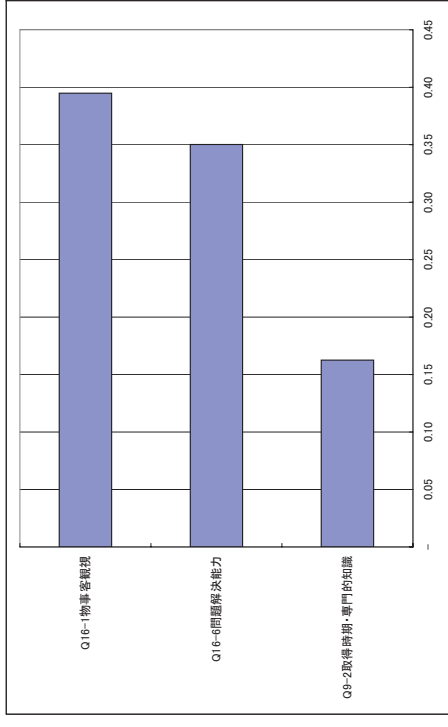
資料 4-3 [重回帰式] 目的変数 仕事スタイル

説明変数名 (自己変数)	偏回帰係数	標準偏回帰係数	F値	P値	判定	T値	標準誤差	偏相関	単相関	符号zz
Q9-2取得時期・専門的知識	0.257	0.163	3.484	0.066 [ ]	1.866	0.137	0.209	0.263		
Q16-6問題解決能力	0.355	0.195	5.9E-09 [***]	3.728	0.168	0.439	0.544			
Q16-8問題解決能力	0.626	0.350	13.892	3.7E-04 [***]	3.728	0.168	0.330	0.330		
定数項	-4.713									

[精度]  
決定係数 R<sup>2</sup> = 0.444  
自由度数 422  
重相関係数 R' = 0.666  
自由度数 422  
決定係数 R<sup>2</sup> = 0.882  
自由度 422

[分散分析表]

変動要因	偏差平方和	自由度	F値	P値	判定
全体変動	101.03	79			
回帰による変動	48.14	3	15.045	9.7E-10 [***]	[**]
回帰からの残差変動	38.80	76	0.743		



資料 4-4 [重回帰式] 目的変数 仕事スタンス

説明変数名 (自己変数)	偏回帰係数	標準偏回帰係数	F値	P値	判定	T値	標準誤差	偏相関	単相関	符号zz
Q9-1取得時期・社会系志願	0.985	0.270	9.270	0.003 [**]	3.045	0.127	0.330	0.436		
Q9-5取得時期・論理的思考	0.272	0.205	1.826	0.181 [ ]	1.344	0.149	0.165	0.546		
Q16-5本質見抜く能力	0.892	0.460	18.765	2.9E-03 [***]	4.44	0.165	0.464	0.546		
定数項	-3.281									

[精度]  
決定係数 R<sup>2</sup> = 0.464  
自由度数 443  
重相関係数 R' = 0.681  
自由度数 443  
決定係数 R<sup>2</sup> = 0.847  
自由度 443

[分散分析表]

変動要因	偏差平方和	自由度	F値	P値	判定
全体変動	101.633	79			
回帰による変動	47.163	3	15.721	2.9E-10 [***]	[**]
回帰からの残差変動	34.470	76	0.717		

