

応用計量経済学(2)

横浜市立大学商学部教授

松浦 克己

大阪大学国際公共政策研究科助教授

Colin McKenzie

第2章 仮説の検定 I

1 帰無仮説と対立仮説

日本の雇用慣行の中でも年功賃金制度は、その特徴の一つとして有名である。年功賃金制度は簡単にいうとある一定の年齢までは年齢と共に賃金が上昇するということを指している。言い換えれば年齢が賃金に影響していることを仮定しているものと言えよう。実際50歳半ば前後まで賃金が上昇し、50代後半以降から賃金が低下するという労働市場に関するグラフを、読者は見られたことがあるだろう。多くのケースで賃金と年齢が関係しているというような推測 (conjecture) を仮説 (Hypothesis) といい、それを統計的に確かめることを仮説の検定 (Hypothesis Testing) という。たとえば以下のようなモデルを考えたとする (なお本章でも仮定A.1からA.5は充たされているものとする¹⁾)。

$$y_i = a + b_1x_{i1} + b_2x_{i2} + \dots + b_jx_{ij} + \dots + b_kx_{ki} + e_i \quad (2.1)$$

x_{ji} が y_i に影響していないならば、 $b_j = 0$ である。逆に x_{ji} が y_i に影響しているならば $b_j \neq 0$ である。 $b_j = 0$ が何%の確率で発生するかを確かめようと

というのが仮説の検定である。検定したい仮説を帰無仮説 (null hypothesis) といい、帰無仮説と論理的に対立し相容れない仮説を対立仮説 (alternative hypothesis) という (帰無仮説を H_0 、対立仮説を H_1 で表すことが多い)。したがって帰無仮説が間違いとされたときは対立仮説が採択される。

この例では、帰無仮説と対立仮説は

$$H_0 : b_j = 0$$

$$H_1 : b_j \neq 0$$

のように表記する。

帰無仮説は検定したい仮説であるが、我々はこのモデルで x_{ji} は y_i に影響しているということを主張したい (だからモデルに必要な変数として x_{ji} が加えられている)。言い換えれば $b_j = 0$ を否定したいと考えている。そこで多くの場合否定したい仮説 (無かったことに帰着させたい) を帰無仮説に選ぶことになる。前回の章で t 検定を紹介したが、 t 検定は計量分析で最も頻繁に行われる仮説の検定である。前回の例で消費は所得の関数であるというとき

$$\text{yearcons}_i = a + b \cdot \text{disposal}_i + e_i$$

のモデルで $b \neq 0$ と主張したいのだから、 $b = 0$ を否定したかったのである。

1) A.1 誤差項の期待値は0である。 $E(e_i) = 0$ for all i
 A.2 誤差項の分散は一定である。 $\text{Var}(e_i) = E(e_i^2) = \sigma^2$ for all i
 A.3 誤差項間に系列相関はない。 $E(e_i e_j) = \text{Cov}(e_i, e_j) = 0$ for all $i \neq j$
 A.4 誤差項は正規分布 (normal distribution) に従う。
 A.5 説明変数はある特定の値を取る非確率変数である。

2 有意水準と棄却域

(有意水準)

$b = 0$ であれば、観測された \hat{b} より大きい推定値の絶対値がどの程度の確率(%)で発生するかということが問題となる。

$b = 0$ であれば、観測された \hat{b} より大きい推定値の絶対値が5%の確率で発生するとき有意水準(significance level)は5%であるという。1%の確率で発生するときは有意水準は1%、あるいは10%の確率で発生するときは有意水準は10%であるという。この発生する確率をp値(p-value)という。EviewsはOLSの推計の場合このp値を自動的に出力する。ある特定の1変数の値が、 $b = 0$ であるかどうかはt分布をもとにした検定を行う。このt値のように検定に必要な統計量を検定統計量(test statistic)という。またある有意水準を定めたときそれに対応する統計量を臨界値(critical value)という。

A.1~A.5の仮定が充たされているとき、2.1)式のOLS推定量 b_j を用いると \hat{b}_j と真の値 b_j との差を測るt値は

$$t = \frac{\hat{b}_j - b_j}{\text{se}(\hat{b}_j)} \quad 2.2a)$$

で計算された。これが自由度($n - (k + 1)$)のt分布に従う。なお n は2.1)式の標本数、 $k + 1$ は2.1)式の係数の数である。真の値が0であるとすると、

$$t = \frac{\hat{b}_j}{\text{se}(\hat{b}_j)} \quad 2.2b)$$

を自由度($n - (k + 1)$)のt検定を行うことになる。2.2b)式を利用し、 \hat{b}_j が0と異なるかどうかを検定する。

$H_0: b_j = 0$ の帰無仮説が正しい場合も \hat{b}_j は線形不偏推定量であるから、 \hat{b}_j は0に近く、 \hat{b}_j と比べて $\text{se}(\hat{b}_j)$ は大きくなるのでt値も0に近くなる。 $H_1: b_j \neq 0$ の対立仮説が正しければ、 \hat{b}_j は0から

離れ、 \hat{b}_j に比べて $\text{se}(\hat{b}_j)$ は小さくなるので、t値は0から離れるであろう。ある有意水準を定めた場合、 $|t\text{値}| > \text{臨界値}$ であれば $H_0: b = 0$ ということは、有意水準以下の確率でしか発生しないことになる。このときは $H_1: b \neq 0$ と考えることができそうである。

たとえば前回の消費関数では

$$\text{yearcons} = 253.2427 + 0.1642 * \text{disposal} + \hat{e}$$

2.3)

(13.434) (7.315) () 内はt値

[0.0000] [0.0000] [] 内はp値

(18.850) (0.0224) () 内は標

準偏差

という結果が得られていた。自由度121(サンプルは123、推定される係数は定数項と説明変数の2個である)の有意水準5%のt値は1.98である(したがってこの場合は1.98が臨界値ということになる)。ここで $7.31 > 1.98$ であるから、 $b = 0$ ということは5%以下の確率でしか発生しないことが分かる。このような場合その説明変数は5%水準で統計的に有意(statistically significant)、あるいは5%水準で有意に0と異なる(significantly different from 0)という。disposalのp値が $b = 0$ かどうかを判断するための検定統計量である。disposalのp値は0.000であるから $b = 0$ ということは千回に1回も起きないということが分かる。

推計された \hat{b} が0以外のある特定の値と有意に異なるかどうかは2.2a)式により検定を行うことが可能である。たとえば \hat{b} が0.3と異なるかどうかを見るためには、2.3)式の例では($H_0: b = 0.3$ $H_1: b \neq 0.3$)

$$t = \frac{0.1642 - 0.300}{0.0224} = -6.05 \text{ である。} \quad 2.4)$$

自由度は $b = 0$ のケースと同じなので5%の有意水準で検定すると適切な臨界値は1.98となる。 $|-6.05| > 1.98$ であるから、5%有意水準で \hat{b}

は0.3と有意に異なることになる（この問題については区間推定で後述する）。したがって $H_0: b = 0.3$ は棄却される²⁾。

（棄却域）

ところで我々は検定を行う際はあらかじめ何%の有意水準で判断するかを定める。例えば5%を有意水準と定めた場合、 $H_0: b = 0$ が否定できるとき、その否定できる5%に入る領域を棄却域（rejection region）という。このケースではdisposalの係数が棄却域に入る確率は万に一つということになる。臨界値、有意水準と棄却域の関係を図示したものが図2.1である。

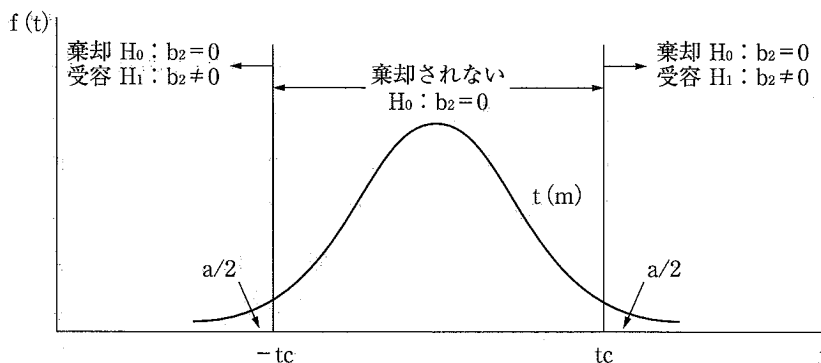
帰無仮説の $|t\text{値}| > \text{臨界値}$ であれば、帰無仮説を誤りと判断しその仮説を棄却する（reject）という。 $|t\text{値}| \leq \text{臨界値}$ であれば、帰無仮説を積極的に否定する根拠がない（その仮説が間違っているという根拠はない）という意味で帰無仮説は棄却されない（not to reject）、あるいは仮説を受容する（accept）という。すなわち仮説を棄却するということは、それを間違いと判断することである。他方仮説を棄却できないということは、必ずしもそれが間違いとは言えないという意味であり、その仮説を積極的に正しいと認めているわけではない。その意味で「棄却する」というのは「棄却できない」ということより強い概念である。

3 片側検定と両側検定

$H_0: b = 0$ という帰無仮説に対する対立仮説を考える場合、それと論理的に対立する仮説としては、 $H_1: b \neq 0$ の他に $H_1: b > 0$ と $H_1: b < 0$ がある。経済理論で消費は所得の関数であると言ったとき、消費の限界性向は正であると考えている。逆に金利が上昇すれば住宅投資は減少すると言ったときは、住宅投資は金利の減少関数であると考えていることになる。このように理論的にある説明変数の係数の符号が正（または負）と予想されることがある。その場合は対立仮説としては $H_1: b > 0$ （符号条件が正と予想されるとき）、あるいは $H_1: b < 0$ （符号条件が負と予想されるとき）を考えてやればよい。

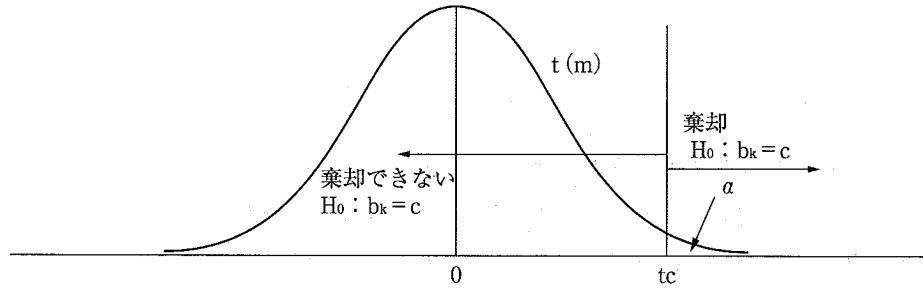
その棄却域は $H_1: b > 0$ であれば上側、 $H_1: b < 0$ であれば下側を用いて行う。なぜならば $b > 0$ （ $b < 0$ ）が正しいとすれば、 b は正（負）で0から離れ、 t 値も同一のはずだからである。棄却域の片方を用いて検定する場合を片側検定（one sided test, one-tailed test）という。棄却域の上側と下側の双方を同時に用いる検定を両側検定（two sided test, two-tailed test）という。 $H_1: b \neq 0$ の対立仮説は、 $b > 0$ と $b < 0$ の二つの概念を包含するから両側検定を行うことになる。

図2.1 臨界値、有意水準、棄却域



2) 帰無仮説が正しければ、(0.1642-0.30)の値は有意に0と異なるはずである。

図2.2 片側検定の臨界値、有意水準、棄却域



たとえば正が予想される場合の有意水準5%の棄却域は図2.2で示される。

先の消費関数の例に戻ろう。帰無仮説は $H_0: b = 0$ 、対立仮説は $H_1: b > 0$ である。t値の導出は前と同様である。有意水準を5%と定めた場合、棄却域は上側のみを見てやればよい。t分布の自由度121の上側5%有意水準は1.66である。7.31 > 1.66であるから $H_0: b = 0$ は棄却され、 $b > 0$ という対立仮説が採択される。なおt分布は左右対称であるからp値は両側検定の1/2となる(Eviewsで表示されるp値は両側検定を行ったときのp値が表示される。それを1/2倍してやればよい)。

4 第1種の誤りと第2種の誤り

帰無仮説は真か偽のいずれかである。また検定の結果は棄却するか棄却しない(受容する)かのいずれかである。我々が行う検定は、正しい場合もあれば誤ることもある。正しいケースは

- ① 帰無仮説が真、検定結果は仮説を棄却しない(受容する)。
 - ② 帰無仮説が偽、検定結果は仮説を棄却する。である。誤るケースは次の通りである。
 - ③ 帰無仮説が真、検定結果は仮説を棄却する。
 - ④ 帰無仮説が偽、検定結果は仮説を棄却しない(受容する)。
- ③のケースを第1種の誤り(type I error)といい、④のケースの誤りを第2種の誤り(type II

error) という。これをまとめると以下のようにある。

	H_0 は真	H_0 は偽
H_0 を棄却 H_0 を受容	第1種の誤り 正しい検定	正しい検定 第2種の誤り

第1種の誤りを起こす確率は有意水準に等しい。したがって第1種の誤りを起こす確率を低くしようと考えるならば、有意水準をたとえば5%から1%に変更してやればよい。t値は

$$t = \frac{\hat{b}}{se(\hat{b})}$$

である。先の消費関数の例(両側検定)

で有意水準5%に対応する臨界値は1.98、有意水準1%に対応する臨界値は2.62である。 $\hat{b} > 1.98 se(\hat{b})$ または $\hat{b} < -1.98 se(\hat{b})$ であれば5%水準で帰無仮説 $H_0: b = 0$ を棄却する($-1.98 se(\hat{b}) < \hat{b} < 1.98 se(\hat{b})$ であれば帰無仮説を受容する)。1%水準では $\hat{b} > 2.62 se(\hat{b})$ または $\hat{b} < -2.62 se(\hat{b})$ であれば帰無仮説を棄却する($-2.62 se(\hat{b}) < \hat{b} < 2.62 se(\hat{b})$ であれば受容する)。

仮に真の値が $b > 0$ であったとする。その分布は次のように描かれる。有意水準を5%から1%に変更すると $H_0: b = 0$ の帰無仮説を受容してしまう確率(第2種の誤りを起こす確率)が高くなる。つまり第1種の誤りを起こす確率を低下させると、第2種の誤りを起こす確率が高くなる。これから第一種の誤りを起こす確率と第二種の誤り

図2.3.1 type I error と type II error

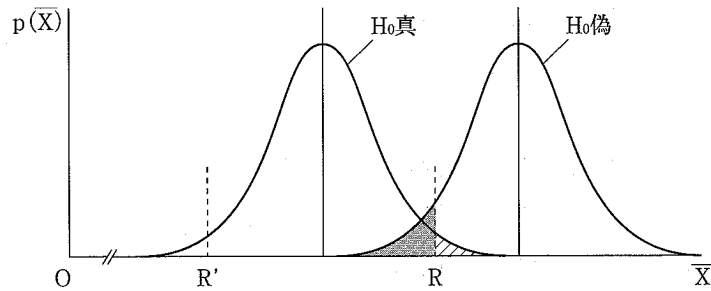
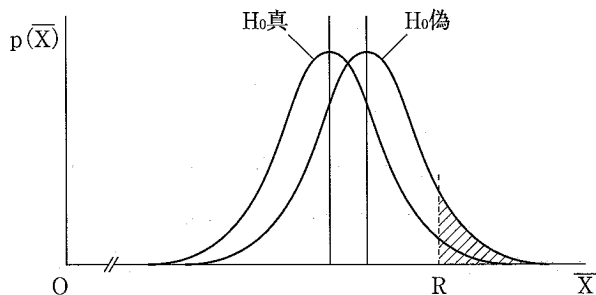


図2.3.2 検出力



を起こす確率は、いずれか一方を低くしようとすると他方が高くなるトレードオフの関係に立つことが分かる (図2.3.1参照)。

そこで通常検定を行う場合には、第1種の誤りを起こす確率を一定として第2種の誤りを起こす確率を低くするようにする。このとき

$$1 - p \text{ (第2種の誤りの確率)} \quad 2.4)$$

を検出力 (power of test) という (図2.3.2参照)。

5 点推定と区間推定、点予測と区間予測

5.1 点推定と区間推定

今までの議論では主にあるパラメータの平均値を求めてきた。この様にあるパラメータについて一つの値を求める推計を点推定 (point estimate) という。しかし推定量 (あるいは推定値) は、確率変数であるから、幅を持って考えることもできる³⁾。このようにある一定の確率でパラメータがどの範囲に存在するかを推計することを、

区間推定 (interval estimate) という。つまり区間推定は、得られた推定値がある確率 (例えば95%) で最大最小でどの値をとるのかということである。それは一般的に次のように求められる。

$$\hat{y}_i = \hat{a} + \hat{b}_1 x_{i1} + \hat{e}_i \quad 2.5)$$

が得られたとする。

ある確率を定める。この定めた確率を信頼度 (level of confidence)、あるいは信頼係数という。このとき2.2a) の結果を利用すると信頼度 (例えば95%) にパラメータが入る区間は

$$P((\hat{b} - b) \pm t^* se(\hat{b})) = 0.95 \quad 2.6)$$

t^* は、(1 - 信頼度) すなわち有意水準に対応するt値の臨界値。

2.6) 式のカッコ内を書き直すと

$$\hat{b} - t^* se(\hat{b}) < b < \hat{b} + t^* se(\hat{b}) \quad 2.7)$$

となる。この区間が信頼区間 (confidence interval) である。先の消費関数例ではbの95%の信頼区間が $0.120 < b < 0.209$ となることを読者は確かめられたい (先に我々はdisposalの点推定で得られた値0.16が0.30と5%水準で有意異なることを見た。これはdisposalの係数の推定値が信頼度95% (5%の有意水準) の信頼区間に0.30を含んでいない例である。0.209 < 0.30を考えれば容易に理解できよう)。

3) 連続確率変数がある特定の一点をとる確率は0である。

5.2 点予測と区間予測

$y_i = a + b_1 x_{1i} + e_i$ を推計し、その結果次の予測が得られたとする。

$$\hat{y}_i = \hat{a} + \hat{b}_1 x_{1i} \quad (2.8)$$

これは a や b_1 の求められた点推定値(平均値)をそのまま利用して、被説明変数の動きを予測しようというものである。これを点予測(point prediction)という。点推定(平均値)を用いた予測であるから、平均での予測ということができる。

x_i の値が x_0 のとき y_0 のモデルが $y_0 = a + b_1 x_{10} + e_0$ とし、 y_0 の予測値を \hat{y}_0 とすると、予測誤差は

$$\hat{e}_0 = \hat{y}_0 - y_0 = (\hat{a} - a) + (\hat{b}_1 - b_1)x_{10} - e_0 \quad (2.9)$$

となる。 e_0 は標準的仮定を充たすとする。また仮定により a と b_1 は不偏推定量であるから $E[(a - \hat{a}) + (\hat{b}_1 - b_1)x_{10}]$ は0となる。したがって2.9)式の期待値は0となるので、 \hat{y}_0 は y_0 の不偏予測推定量である。

予測誤差の分散は

$$V(\hat{y}_0 - y_0) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_{10} - \bar{x}_1)^2}{\sum (x_{1i} - \bar{x}_1)^2} \right] \quad (2.10)$$

で得ることができる。 σ^2 の不偏推定量 s^2 で置き換えると

$$\hat{V}(\hat{y}_0 - y_0) = s^2 \left[1 + \frac{1}{n} + \frac{(x_{10} - \bar{x}_1)^2}{\sum (x_{1i} - \bar{x}_1)^2} \right] \quad (2.11)$$

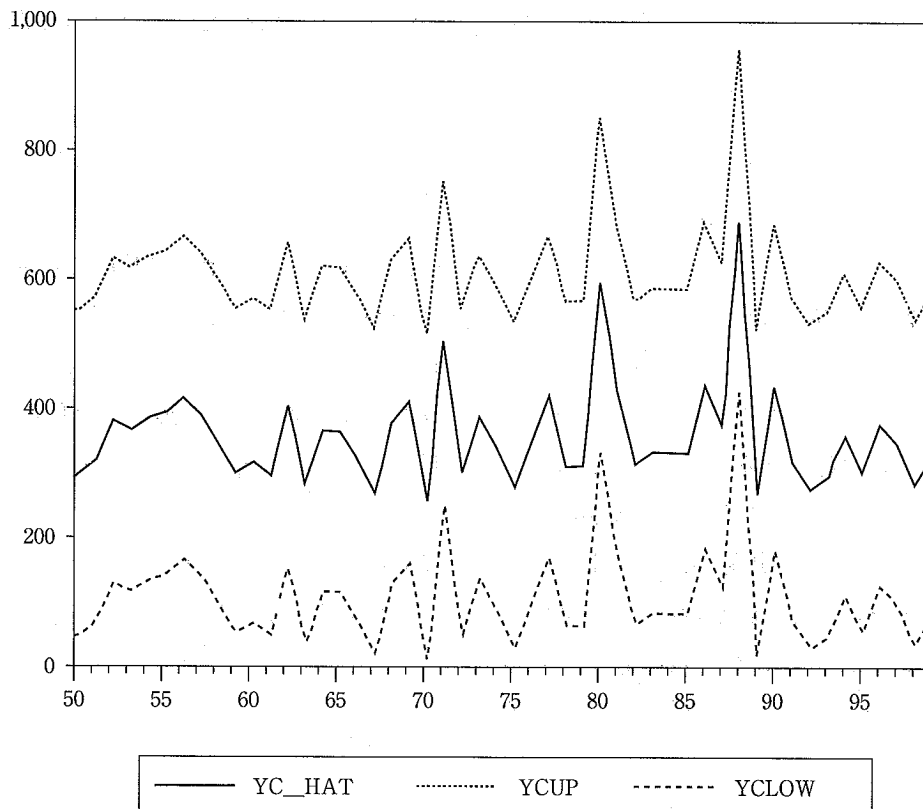
となる。その平方根が $\sqrt{\hat{V}(\hat{y}_0 - y_0)}$ で予測の標準誤差(standard error of forecast)である。これを仮に $se(f)$ と書こう。

$$\frac{(\hat{y}_0 - y_0)}{se(\hat{f})} \quad (2.12)$$

は自由度 $n - 2$ のt分布に従う(この例では推定のために利用した標本数が n で、推定される係数は2個である)。

これから区間推定の場合と同様に信頼度を設定して、その範囲内で予測がどの程度の幅を持つか

図2.4 予測値と信頼区間



という予測の信頼区間を設定できる。この様にして得られた予測を区間予測 (interval prediction) という。

$$P(\hat{y}_0 - t_c s_e(f) \leq y_0 \leq \hat{y}_0 + t_c s_e(f)) = A \quad 2.13$$

ここで t_c は求めようとする信頼度の臨界値、 A は信頼度となる。予測の95%の信頼区間を正確に計算するために t_c を調べる必要があるが、標本数 (n) が十分に大きいときは、 t_c は2前後となる。消費関数の信頼度 ± 2 標準偏差での信頼区間を掲げておく。そのプログラムは以下の通りである (' で始まる行はEviewsでは注釈行で、その行の作業は行わない。プログラムの解説に便利である)。

```
equation eq_1_1.ls yearcons c disposal
' 予測値をyc_hat、その標準誤差をyc_seとする
(' _hat, _seはEviewsの予測値と予測標準誤差を表す既定の表現方法)
```

```
eq_1_1.fit yc_hat yc_se
```

```
'  $\pm \sigma$  を求める。
```

```
genr ycup = yc_hat + 2 * yc_se
```

```
genr yclow = yc_hat - 2 * yc_se
```

```
' 予測値と  $\pm 2 \sigma$  の信頼区間を図示する
```

```
plot yc_hat ycup yclow
```

6 多重共線関係

モデルが複数の説明変数を含む場合、多重共線関係 (multicollinearity) といわれる困難な問題を生じることがある。たとえば

$$y_i = a + b_1 x_{1i} + b_2 x_{2i} + e_i \quad 2.14$$

$$x_{2i} = dx_i \quad 2.15$$

のように x_{2i} と x_{1i} が線形従属である (完全に相関している) としよう。2.14) 式を2.15) 式に代入すると

$$y_i = a + b_1 x_{1i} + b_2 (dx_{1i}) + e_i$$

$$= a + (b_1 + b_2 d) x_{1i} + e_i \quad 2.16$$

となり b_1 と b_2 を推計することはできない (第1章で多重回帰の係数を求めた1.35) 式の分母 $S_{11}S_{22} - S_{12}^2 = 0$ となることから容易に確かめることができる⁴⁾。この場合を完全な多重共線関係 (perfect multicollinearity) という。このように説明変数間に完全に相関があるのは希であろうが、相関が高い場合も類似の問題が起こりうる。多重共線関係が存在する時以下のような問題が生じる。

- ① サンプルを増減すると推定値が大きく変化する。
- ② 説明変数を入れ替えると推定値が大きく異なる (時には正負の符号が逆転することがある)。
- ③ 推定値の分散 (標準誤差) が大きくなり、本来統計的に有意な変数を非有意と誤ることがある。

このために多重共線関係にあるときは、その多重共線関係にある各説明変数の影響を個別には捉えることができなくなる (多重共線関係にないその他の説明変数は影響を受けない)。

さらに多重共線関係は説明変数間の相関が高くない場合にも起こりうる。

より一般的に、 $y_i = a + b_1 x_{1i} + b_2 x_{2i} + \dots + b_k x_{ki} + e_i$ のケースで j 番目のパラメータの分散 $V(b_j)$ を考える。 σ^2 を誤差項の分散、 $S_{jj} = \sum_i (x_{ji} - \bar{x}_j)^2$ 、 R_{jj}^2 を x_{ji} を他の全ての説明変数に回帰したときの R^2 とする。以下の結果を得る。

$$V(\hat{b}_j) = \frac{\sigma^2}{S_{jj}(1 - R_{jj}^2)} \quad 2.17$$

これから \hat{b}_j の分散は σ^2 が大きいほど、 S_{jj} が小さいほど、また R_{jj}^2 が高いほど大きくなる事が分かる。したがってたとえ説明変数間の相関が低くとも、 σ^2 や S_{jj} の値によっては推定された分散は

4) 多重回帰の前提として $S_{11}S_{22} - S_{12}^2 \neq 0$ 、あるいは完全な多重共線関係のないことを仮定する研究者もいる。

大きくなり、多重共線関係は起こりうる（多重共線関係が起きたと判断したときEviewsは、“Near singular matrix” というエラーメッセージを出す）。

多重共線関係の尺度としていくつかの指標が提案されているが、一致した見解は得られていないのが実状である。また多重共線関係問題の解決には、（可能であれば）サンプルを増やすこと（ S_{ij} が大きくなる）、あるいは不必要な変数を落とすことなど様々な案が出されているが、決め手はないのが現状である。

（どのように報告するか）

しかしRule of Thumbとして次のように報告することは有益である。 x_1 と x_2 が多重共線関係にあるとしても、その全ての説明変数を含む推計、多重共線関係にある（と疑われる）各1個の変数を落とした推計を報告することである。この例を後ほど実際に見てみることにする。また多重共線関係にあるときも、点予測については全ての説明変数を含む回帰の場合、多重共線関係の影響は受けないので、点予測に用いることはできることが知

られている。

たとえば金融資産を入れた消費関数についてみよう。disposal、money、numberの相関係数を求めよう。そのためのコマンドは次の通りである。

`cor (p) disposal money number`

そうするとdisposalとmoneyの相関係数は0.43であり、それほど高いわけではないということを読者は確認されたい⁵⁾。

次に前章のequation 1_3からdisposalを除いて推計してみよう。

`equation eq 1_4 .ls yearcons c number money`

この結果は表2.1に示すとおりである。moneyの係数は0.02となり、t値は4.01、p値は0.0001と変わる。金融資産は1%水準でも有意に正の影響を消費に与えている。このケースはdisposalとmoneyで多重共線関係が起きていた可能性があることを示唆している（断言はできないが）。

この様に多重共線関係が疑われるとき、個々の説明変数の影響が明確でなくなるので、我々は判断に迷うことが多い。しかしRule of Thumbに従い、このケースでは次の3通りを報告することが

表2.1 金融資産を入れ所得を除いた推計例

Dependent Variable: YEARCONS				
Method: Least Squares				
Sample: 1 123				
Included observations 123				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	207.2845	28.32099	7.319113	0.0000
NUMBER	40.07099	7.993575	5.012899	0.0000
MONEY	0.021097	0.005259	4.011341	0.0001
R-squared	0.258245	Mean dependent var		363.8049
Adjusted R-squared	0.245882	S.D. dependent var		149.4013
S.E. of regression	129.7400	Akaike info criterion		12.59303
Sum squared resid	2019896.	Schwarz criterion		12.66162
Log likelihood	-771.4713	F-statistic		20.88923
Durbin-Watson stat	2.074385	Prob(F-statistic)		0.000000

5) 共分散はcov (p) disposal money numberで求めることができる。

有益である。

$$\begin{aligned} \text{yearcons} &= 179.28 + 0.123 \text{disposal} \\ &\quad (6.76) \quad (4.96) \\ &+ 28.362 \text{number} + 0.009 \text{money} + \hat{e}_1 \\ &\quad (3.69) \quad (1.71) \end{aligned}$$

$$\text{AdjR}^2 = 0.370 \quad \text{SER} = 118.61$$

$$\begin{aligned} \text{yearcons} &= 186.23 + 0.142 \text{disposal} \\ &\quad (7.06) \quad (6.34) \\ &+ 26.71 \text{number} + \hat{e}_2 \\ &\quad (3.48) \end{aligned}$$

$$\text{AdjR}^2 = 0.360 \quad \text{SER} = 119.56$$

$$\begin{aligned} \text{yearcons} &= 207.28 + 40.07 \text{number} \\ &\quad (7.32) \quad (5.01) \\ &+ 0.021 \text{money} + \hat{e}_3 \\ &\quad (4.01) \end{aligned}$$

$$\text{AdjR}^2 = 0.246 \quad \text{SER} = 129.74$$

このとき点予測については全ての説明変数を含む回帰の場合、多重共線関係の影響は受けないので、点予測に用いることはできることが知られていることは前に述べたとおりである。

7 ダミー変数

ある条件を満たす場合を1、ある条件を満たさない場合は0となるような変数をダミー変数 (Dummy Variable) という。たとえば女性であれば1、男性であれば0というのはその一例である。更にあるグループを属性により3個以上に分けてダミーを作ることもできる。たとえば勤務先の従業員数により99人以下、100~499人、500人以上と区分すれば、3個のダミー変数を作ることができる。

(定数項ダミー)

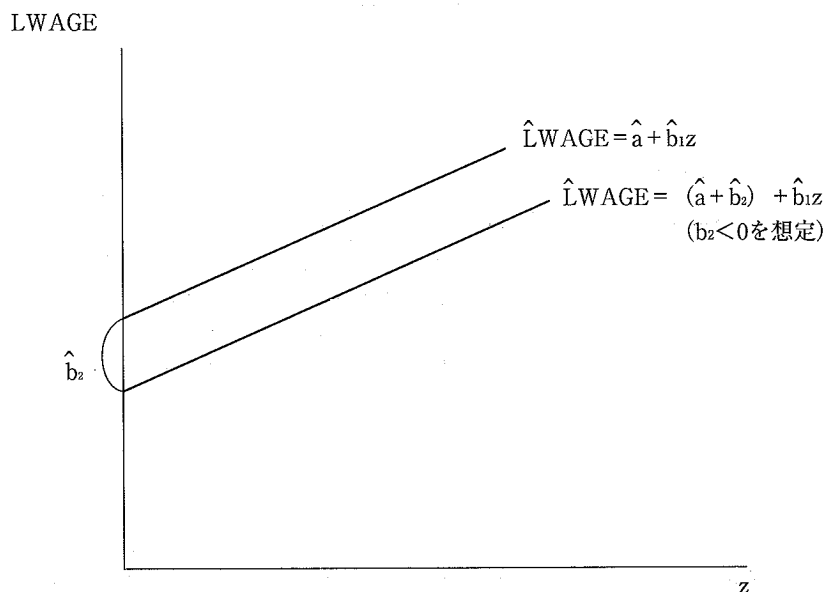
日本の労働市場では男女間の賃金格差があることが指摘されている。これは他の条件を一定にして女性の賃金が男性よりも低いことを意味している。具体的には以下のようなモデルで

$$\text{lwage}_i = a + b_1 z_i + b_2 \text{female}_i + e_i \quad (2.18)$$

lwageは賃金の対数、 z は女性以外の説明変数、femaleは女性の場合=1、男性の場合=0となるダミー変数。 a 、 b_1 、 b_2 は推計すべきパラメータ。 e_i は誤差項

男女間の賃金格差仮説は $b_2 < 0$ を想定していることになる。女性と男性についてダミー変数を用

図2.5 定数項ダミーの効果



いて書き分けると次のようになる。

$$lwage_i = (a + b_2) + b_1z_i + e_i \quad \text{if female} = 1 \quad (2.19)$$

$$lwage_i = a + b_1z_i + e_i \quad \text{if female} = 0 \quad (2.20)$$

2.19) 式と2.20) 式を比較すると定数項の部分
が異なることが分かる。他の説明変数に係る b_1 は
共通であるから、両式は b_2 の部分だけ $lwage$ の予
測値が異なるので、回帰直線は b_2 だけ平行にシフ
トする (図2.5 参照)。

このように定数項の部分だけに影響するダミー
変数を定数項ダミー (intercept dummy variable)
という。なお男性 (female = 0) と比べて女性の
賃金はどうなっているかを見るものであるので、
ダミー変数が0となるものを基準値あるいは既定
値 (default) ということがある。

これをEviewsで実際に見てみよう。データは
「家計の金融資産選択 (1996年)」である。

- 1 workfile a: labor u 1109
- 2 smpl 1-1109
- 3 read a: labor. dat lwage kigyuu edu age
female
- 4 group group21 lwage age female
- 5 group21. stats
- 6 series age 2 = age ^ 2
- 7 equation eq 2_1. ls lwage c age age 2

female

1行目で作業領域を設定している。2行目でサ
ンプル数が1109であることを指定している。その
変数が $lwage$ (世帯主の収入、万円の対数値)、
 $kigyuu$ (勤め先の従業員数ダミー、1 = 4人以
下、2 = 5—29人、3 = 30—99人、4 = 100—499
人 5 = 500人以上)、 edu (教育歴、1 = 中卒、2
= 高卒、3 = 短大卒、4 = 大卒、6 = その他)、
 age (年齢)、 $female$ (世帯主が女性 = 1、男性 =
0) であることを指定している。

7行目で $lwage_i = a + b_1age_i + b_2age_2_i + b_3female_i + e_i$ を推計するコマンドを指示している (方程式
の名前は $eq\ 2_1$ である)。

結果は表2.2に示す通りである。 $female$ の係数
は-0.39である。 t 値は-7.59であるから1%水
準で統計的に有意である。 p 値は0.0000であるか
ら $b_3 = 0$ の帰無仮説は、両側検定でも片側検定で
も、強く棄却されている。この結果では女性の賃
金は男性に比べて0.39 (対数値) 低くなっており、
男女間の賃金格差仮説は支持される。

(複数のダミー変数)

定数項ダミーを2個以上作ることも可能である。
企業規模による賃金格差がいわれている。これは
大企業、中堅企業、中小企業の従業員の間で他の

表2.2 定数項ダミーの例

Dependent Variable: LWAGE				
Method: Least Squares				
Included observations: 1109				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	4.398062	0.186177	23.62299	0.0000
AGE	0.078317	0.008644	9.060278	0.0000
AGE 2	-0.000745	9.64E-05	-7.729260	0.0000
FEMALE	-0.391149	0.051521	-7.592103	0.0000
R-squared	0.188008	Mean dependent var		6.275685
Adjusted R-squared	0.185804	S.D. dependent var		0.497480
S.E. of regression	0.448890	Akaike info criterion		1.239524
Sum squared resid	222.6602	Schwarz criterion		1.257598
Log likelihood	-683.3159	F-statistic		85.28362
Durbin-Watson stat	1.642526	Prob(F-statistic)		0.000000

条件を一定として賃金水準が異なることを意味している。このようなケースでは大企業、中堅企業と中小企業の3個のグループ間で違いがあるかどうかを課題とする。それをダミー変数を用いることにより検証してみよう。

まず企業間の賃金格差をEviewsで推計してみよう。なお説明の便宜のためにここでも注釈行をつける。

前回の論理演算子の所で説明したように、

series A変数 = B変数 = ある条件

とすれば、B変数が条件を満たすとき、 $A = 1$ (満たさないとき $A = 0$) という変数が作成された。これを利用する。企業規模は3個に区分され、大企業、中堅企業と中小企業はそれぞれ500人以上、100-499人と99人以下の企業とする。99人以下の企業を基準値として推計してみよう。

```
8 ' 1—4人企業ダミー
9 series emp1 = kigyou = 1
10 ' 5—29人企業ダミー
11 series emp5 = kigyou = 2
12 ' 30—99人企業ダミー
```

```
13 series emp30 = kigyou = 3
14 ' 1—99人の企業ダミー
15 series emp13 = kigyou <= 3
16 ' 100—499人企業ダミー
17 series emp100 = kigyou = 4
18 ' 500人以上企業ダミー
19 series emp500 = kigyou = 5
20 ' 記述統計量
21 group group22 emp1 emp5 emp30
   emp13 emp100 emp500
22 group22. stats
23 ' lwagei = a + b1agei + b2age 2i + b3femalei +
   b4emp500i + b5emp100i + eiを推計するコマ
   ンド
24 equation eq2_2: ls lwage c age age 2
   female emp500 emp100
```

24行は、(2.18)式にならうと以下のように書くことができる。

$$lwage_i = a + b_1age_i + b_2age_{2i} + b_3female_i + b_4emp500_i + b_5emp100_i + e_i \quad (2.21)$$

従ってダミー変数の効果は以下のように表すこと

表2.3 複数のダミーの例

Dependent Variable: LWAGE
Method: Least Squares
Sample: 1 1109
Included observations: 1109

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	4.407635	0.175930	25.05334	0.0000
AGE	0.068743	0.008221	8.362310	0.0000
AGE 2	-0.000622	9.18E-05	-6.771728	0.0000
FEMALE	-0.322238	0.049057	-6.568574	0.0000
EMP100	0.151617	0.033619	4.509907	0.0000
EMP500	0.343299	0.029595	11.59970	0.0000
R-squared	0.276298	Mean dependent var		6.275685
Adjusted R-squared	0.273017	S.D. dependent var		0.497480
S.E. of regression	0.424168	Akaike info criterion		1.128020
Sum squared resid	198.4499	Schwarz criterion		1.155132
Log likelihood	-619.4872	F-statistic		84.22147
Durbin-Watson stat	1.657493	Prob(F-statistic)		0.000000

ができる。

$$\begin{aligned} \text{lwage}_i &= (a + b_4) + b_2 z_i + e_i \text{ emp500} \\ &= 1, \text{ emp100} = 0 \text{ のケース} \end{aligned} \quad 2.22)$$

$$\begin{aligned} \text{lwage}_i &= (a + b_5) + b_2 z_i + e_i \text{ emp100} \\ &= 1, \text{ emp500} = 0 \text{ のケース} \end{aligned} \quad 2.23)$$

$$\begin{aligned} \text{lwage}_i &= a + b_2 z_i + e_i \text{ emp500} = 0, \text{ emp100} \\ &= 0 \text{ のケース} \end{aligned} \quad 2.24)$$

このケースでは99人以下の企業が基準となっているので、 b_4 は500人以上企業と99人以下企業の賃金格差、 b_5 は100-499人企業と99人以下企業の勤労者の賃金格差を表している。

なおここでは99人以下の企業ダミー (emp13) が取り上げられていない。これはこの変数も取り入れた場合、各経済主体に取り

$$1 = \text{emp500}_i + \text{emp100}_i + \text{emp13}_i \quad 2.25)$$

となるので、完全な多重共線関係となり推計できないからである。

推計結果は表2.3に示す通りである。emp100とemp500の係数は正である。しかもt値は4.51と11.60でありいずれも1%水準で有意である。100-499人、あるいは500人以上の企業に勤務する勤労者の賃金は、99人以下の企業に勤務する勤労者の賃金よりそれぞれ0.15, 0.34 (対数値) 高いことが分かる。

(交差項、係数ダミー)

ダミー変数とダミー変数をかけた変数を考えることもできる。あるいはダミー変数と他の連続変数をかけることもできる。複数の変数をかけあわせた変数を交差項 (interaction variable) という。後者を特に係数ダミー (slope dummy variable) ということがある。

たとえば女性ダミー*500人以上勤務ダミーとすれば、それは女性でかつ500人以上の企業に勤務する人を指すダミーである (第1章のoldmanを思い出してほしい)。

係数ダミーはそれに係る連続説明変数の傾きを

変えることになる。たとえば年功賃金といっても女性の年功カーブと男性の年功カーブでは異なるかもしれない (男女間の賃金格差は年齢と共に拡大するという事柄もしばしば言われている)。そこで次のような関数を考える。

$$\begin{aligned} \text{lwage}_i &= a + b_1 \text{age}_i + b_2 \text{age } 2_i + c_1 (\text{female}_i * \text{age}_i) \\ &\quad + c_2 (\text{female}_i * \text{age } 2_i) + d_2 z_i + e_i \end{aligned} \quad 2.26)$$

男性であればfemale = 0 であるから

$$\text{lwage}_i = a + b_1 \text{age}_i + b_2 \text{age } 2_i + d_2 z_i + e_i \quad 2.27)$$

となる。女性であればfemale = 1 であるから

$$\begin{aligned} \text{lwage}_i &= a + (b_1 + c_1) \text{age}_i + (b_2 + c_2) \text{age } 2_i \\ &\quad + d_2 z_i + e_i \end{aligned} \quad 2.28)$$

となる。女性の年齢の効果は、ageについては $(b_1 + c_1)$ 、age 2については $(b_2 + c_2)$ となるので、その傾きは男性と比べてダミー変数の推定値である c_1 、 c_2 だけ異なることになる。

8 F検定と線形制約

8.1 0制約の例

年功序列賃金制度が存在しないならば、年齢は賃金に統計的に有意な影響を与えていないであろう。言い換えれば以下の式で

$$\begin{aligned} \text{lwage}_i &= c + a_1 \text{age}_i + a_2 \text{age } 2_i + b_1 \text{emp500}_i \\ &\quad + b_2 \text{emp100}_i + b_3 \text{female}_i + e_i \end{aligned} \quad 2.29)$$

$$a_1 = a_2 = 0 \quad 2.30)$$

の制約が成立しているであろう ($a_1 = a_2 = 0$ のようにある係数を0と置くことを0制約という)。つまり次の帰無仮説が成立しているはずである (このように2個以上の仮説を同時に検定することを複合仮説検定 (joint hypothesis test) という)。

$$H_0 : a_1 = a_2 = 0 \quad 2.31)$$

対立仮説は次のようである。

$$H_1 : H_0 \text{ ではない} \quad 2.32)$$

(具体的には $a_1 \neq 0$ または $a_2 \neq 0$ の少なくとも一方が成立する)。

この複合仮説を検定するとき、我々はF検定を

行う。帰無仮説が正しければ、次式が成立しているであろう。

$$lwage_i = c + b_1 emp500_i + b_2 emp100_i + b_3 female_i + e_i \quad (2.33)$$

制約のあるモデル (restricted model, ここでは2.33式) の残差平方和を RSS_R とする。制約のないモデル (unrestricted model, ここでは2.29式) の残差平方和を RSS_U とする。

$$RSS_R \geq RSS_U \quad (2.34)$$

が必ず成立する (age と age 2 がいずれも lwage と完全に無相関であれば等号が成立し、何らかの相関があれば不等号となる)。したがって帰無仮説が正しければ ($RSS_R - RSS_U$) の値は小さくなるであろう。逆に帰無仮説が間違っている (制約が有効ではない) ときはこの値は大きくなるであろう。このとき次のF検定統計量 (F test statistic) が導かれる。

$$F = \frac{(RSS_R - RSS_U) / r}{RSS_U / (n - k)} \quad (2.35)$$

この値は自由度 r 、 $(n - k)$ のF分布に従うことが知られている ($F_{r, n-k}$ と書く、2つの自由度があるという点で正規分布やt分布とは異なる)。なお r は帰無仮説にある制約の数、 n はサンプル数、 k は制約のないモデルの係数の数。2.34) より $F \geq 0$ である。帰無仮説が正しければFの値は小さくなるであろう。逆に帰無仮説が間違っているときはF値は大きくなるであろう。

2.29) 式と2.33) 式を各々OLSで推計し、その残差を用いて2.35) 式のF値を計算することができる。なたこれをEViewsはコマンドで行うことができる。読者は24行の結果を表示してほしい。View/Coefficient Tests/Redundant Variable-Likelihood Ratioを選択し、画面に制約条件である2個の説明変数age age 2を記入し実行してほしい (図2.6、2.7参照)。

以下の結果を得るだろう (表2.4参照)

F値は95.6である。有意水準5% (1%) で自由度2, 1103のF統計量は約3.00 (約4.61) である

図2.6

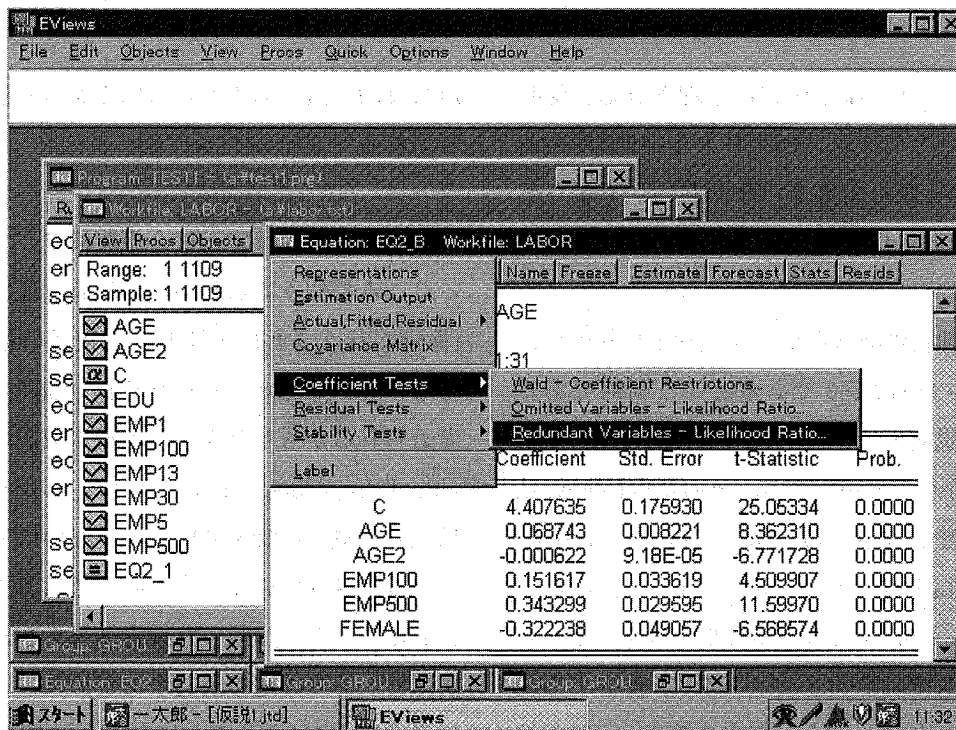


図2.7 画面選択の例

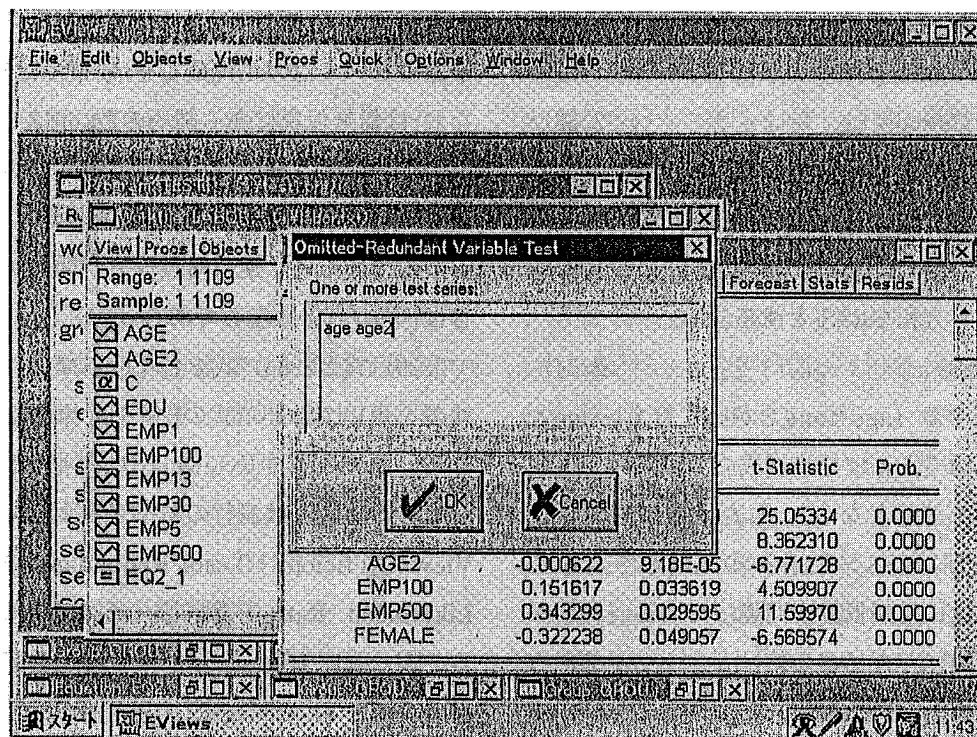


表2.4 0制約のF検定の例 (ageの係数=age2の係数=0)

Redundant Variables: AGE AGE 2

F-statistic	95.59136	Probability	0.000000
Log likelihood ratio	177.2688	Probability	0.000000

から、 $a_1 = a_2 = 0$ の帰無仮説は棄却される。帰無仮説が正しいとして、95.6より大きいF値が起るp値は0.0000となっている（なお画面に表示されるLog Likelihood Ratioの意味については後述する）。

なおゼロ制約の極端な場合が定数項以外の全ての説明変数の係数が0であるという制約である（モデルが説明力を全く持たないケース）。これもF検定で行うことができる。その値はEViewsは自動的に計算し表示する（結果下欄のF-statistic Prob（F-statistic）がこれに該当する。2.29）式において $a_1 = a_2 = b_1 = b_2 = b_3 = 0$ を検定すると表2.3のケースではF-statistic 84.2が該当する統計量である。有意水準5%で自由度5, 1103のF統

計量は約2.21であるから、明らかに帰無仮説は棄却される）。

8.2 単独仮説検定と複合仮説検定

ところで個々の係数の統計的有意度についてはt検定で行うことは説明した。複数の検定についても、 $a_1 = 0$, $a_2 = 0$ とそれぞれ単独に行うことが考えられるかもしれない（単独の仮説検定をindividual hypothesis testという）。しかしこの場合は他の変数の影響を考慮していないので、その効果が排除されていない。そのために多重共線関係にあるときなどは誤った検定結果をもたらすことがある。その問題を回避するために、複合仮説の検定はF検定によることになる。

表2.5 係数ダミーを入れた推計例と複合仮説検定、多重共線問題

Dependent Variable: LWAGE

Method: Least Squares

Sample: 1 1109

Included observations: 1109

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	4.267078	0.189909	22.46904	0.0000
AGE	0.074418	0.008836	8.422197	0.0000
AGE 2	-0.000675	9.81E-05	-6.878043	0.0000
FAGE	-0.022563	0.025540	-0.883435	0.3772
FAGE 2	0.000149	0.000301	0.495768	0.6202
FEMALE	0.324516	0.506588	0.640592	0.5219
EMP100	0.151536	0.033578	4.512936	0.0000
EMP500	0.341322	0.029568	11.54375	0.0000
R-squared	0.281005	Mean dependent var		6.275685
Adjusted R-squared	0.276434	S.D. dependent var		0.497480
S.E. of regression	0.423170	Akaike info criterion		1.125101
Sum squared resid	197.1589	Schwarz criterion		1.161250
Log likelihood	-615.8683	F-statistic		61.47215
Durbin-Watson stat	1.665773	Prob (F-statistic)		0.000000

表2.6.1 fageの係数=fage 2 の係数=0 のF検定

Redundant Variables: FAGE FAGE 2

F-statistic	3.604506	Probability	0.027522
Log likelihood ratio	7.237723	Probability	0.026813

表2.6.2 fage 2 の係数=0 の検定と他の変数の推計値

Redundant Variables: FAGE 2

F-statistic	0.245786	Probability	0.620157	
Log likelihood ratio	0.247544	Probability	0.618810	
FAGE	-0.010042	0.003804	-2.639696	0.0084
FEMALE	0.086637	0.162438	0.533354	0.5939
	(係数)	(標準偏差)	(t値)	(p値)

表2.6.3 fageの係数=0 の検定と他の変数の推計値

Redundant Variables: FAGE

F-statistic	0.780457	Probability	0.377194	
Log likelihood ratio	0.785850	Probability	0.375358	
FAGE 2	-0.000114	4.49E-05	-2.535712	0.0114
FEMALE	-0.115039	0.095246	-1.207809	0.2274
	(係数)	(標準偏差)	(t値)	(p値)

その例を2.24) 式の女性の年功の推計結果をもとに示すことにしよう。

```

25 ' 係数ダミーの作成
26 series fage = age * female
27 series fage 2 = age 2 * female
28 ' lwagei = a + b1agei + b2age 2i + c1(femalei
agei) + c2(femaleiage 2i) + dzi + eiを推計す
るコマンド
29 equation eq 2_3.ls lwage c age age 2
fage fage 2 female emp100 emp500

```

結果は表2.5に示す通りである。

fageのt値は-0.88, fage 2のt値は0.5である。

いずれも統計的に有意な結果は得られていない(p値は0.38と0.62である)。そこでこの2つの係数が同時にゼロという制約を置こう(age、age 2を除いて推計する)。

先ほど同様に読者は、View/Coefficient Tests/Redundant Variable-Likelihood Ratioを選択し、redundant variableとしてfage fage 2を指定してほしい。表2.6.1の結果が得られるであろう。

F値は3.60であるから5%水準で $c_1=c_2=0$ の帰無仮説は棄却されている(p値は0.028である)。そこでfage, fage 2の各1変数を除いた推計を次に行おう。

再びView/Coefficient Tests/Redundant Variable-Likelihood Ratioを選択し、fage (あるいはfage 2)を指定する(結果は表2.6.2と2.6.3参照)。

一方だけを取り上げるケースではfageは1%水準で、fage 2は5%水準で有意である(しかし表2.2では1%水準で有意であったfemaleはいずれも有意ではない)。female、fage、fage 2の統計的有意水準が説明変数の組み合わせにより大きく変わっている。これは多重共線関係が起きているときの典型的な症例である。

このことは改めて2つの課題を示すものである。

一つは複合仮説検定(F検定)を行う必要があるにもかかわらず、単独仮説検定(t検定)で替えるときは結果を誤る可能性があるということである。一つは多重共線関係にあるとき(あるいは多重共線関係の存在が疑われるとき)は、その変数毎の組み合わせた結果を報告しないと、解釈を誤る可能性があるということである。

この表2.6.1から表2.6.3の結果はfemale、fageとfage 2の間に多重共線関係が起きていることを示すものである。このような場合個々の変数についてt検定で判断を下すことが誤りであることを示す具体例である。

9 不要な変数を入れる場合、必要な変数を落とす場合

真のモデルは必ずしも自明ではない。そこで我々は様々な試行錯誤を行うことになる。その時不要な変数をモデルに入れたり(inclusion of an irrelevant explanatory variable)、逆に必要な変数を落とす(omission of a relevant explanatory variable)ことがある。このようなケースを定式化の誤り(mis-specification, specification error)という。不要な変数を入れる場合を過剰定式化、必要な変数を落とす場合を過小定式化という。この問題について考えてみよう。

$$y_i = a + b_1 x_{1i} + e_{2i} \quad (2.36)$$

$$y_i = a + b_1 x_{1i} + b_2 x_{2i} + e_{1i} \quad (2.37)$$

2.36) 式の b_1 の推定量を \hat{b}_1 とし、2.37) 式の b_1 の推定量を b_1 とする。第1章の1.12) 式と1.36) 式から \hat{b}_1 と b_1 は以下のように求められる。

$$\hat{b}_1 = \frac{\sum (x_{1i} - \bar{x}_1) \sum (y_i - \bar{y})}{\sum (x_{1i} - \bar{x}_1)^2} = \frac{\sum (x_{1i} - \bar{x}_1) y_i}{\sum (x_{1i} - \bar{x}_1)^2} \quad (2.38)$$

$$b_1 = \frac{\sum (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) \sum (x_{1i} - \bar{x}_1)(y_i - \bar{y}) - \sum (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) \sum (x_{2i} - \bar{x}_2)(y_i - \bar{y})}{\sum (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) \sum (x_{2i} - \bar{x}_2)(x_{2i} - \bar{x}_2) - [\sum (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)]^2} \quad (2.39)$$

(過小定式化の問題)

2.37) 式が真のモデルであるにも関わらず、

2.36) 式を推計したとする (過小定式化)。すなわち x_2 が y に有意に影響するにも関わらず ($b_2 \neq 0$)、その変数を落としたとしよう。2.37) 式の y_i を 2.38) 式に代入すると

$$\begin{aligned} \hat{b}_1 &= \frac{\sum (x_{1i} - \bar{x}_1) (a + b_1 x_{1i} + b_2 x_{2i} + e_{1i})}{\sum (x_{1i} - \bar{x}_1)^2} \\ &= \frac{b_1 \sum (x_{1i} - \bar{x}_1)^2 + b_2 \sum (x_{1i} - \bar{x}_1) x_{2i} + \sum (x_{1i} - \bar{x}_1) e_{1i}}{\sum (x_{1i} - \bar{x}_1)^2} \\ &= b_1 + \frac{b_2 \sum (x_{1i} - \bar{x}_1) x_{2i}}{\sum (x_{1i} - \bar{x}_1)^2} + \frac{\sum (x_{1i} - \bar{x}_1) e_{1i}}{\sum (x_{1i} - \bar{x}_1)^2} \quad 2.40) \end{aligned}$$

2.40) 式の右辺第 3 項の期待値は仮定により 0 となる。 $\hat{E}(b_1) = b_1$ となるためには第 2 項の期待値が 0 となる必要がある。すなわちこの条件を充たすのは

$$\begin{aligned} \frac{\sum (x_{1i} - \bar{x}_1) x_{2i}}{\sum (x_{1i} - \bar{x}_1)^2} &= \frac{\sum (x_{1i} - \bar{x}_1) (x_{2i} - \bar{x}_2)}{\sum (x_{1i} - \bar{x}_1)^2} \\ &= \frac{\text{Cov}(x_1, x_2)}{V(x_1)} = 0 \quad 2.41) \end{aligned}$$

のときのみである。2.41) 式から $\text{Cov}(x_1, x_2) = 0$ でなければ、 \hat{b}_1 はバイアスを持ちかつ一致性もない。一般的に言ってモデルにとり重要な説明変数を落としてしまうと、最小二乗法の推定量は不偏性と一致性を持たない。

このバイアスは b_2 の符号と、落とした変数と他の説明変数の共分散 (相関) の正負に依存して定まる。たとえば $b_2 > 0$ かつ相関が正であれば正のバイアスを持ち (過大評価)、また $b_2 < 0$ かつ相関が正であれば負のバイアスを持つ (過小評価)。過少定式化の場合 2.36) 式の誤差項は $e_{2i} = b_2 x_{2i} + e_{1i}$ となる。 e_{1i} が標準的線形回帰モデルの仮定を充たすのであれば、 $E(e_{2i}) = b_2 x_2 \neq 0$ である。すなわち過小定式化ではこれは充たされないで、 e_{2i} は標準的線形回帰モデルの仮定 A 1 を充たさないことになる。

(過剰定式化)

次に不要な変数を含んだ場合を考えてみる。

真のモデルが 2.36) 式であるにもかかわらず、2.37) 式を推計したときの問題である。2.36) 式の y_i を 2.39) 式に代入し整理すると

$$b_1 = b_1 + \frac{+ \sum (x_{2i} - \bar{x}_2) (x_{2i} - \bar{x}_2) \sum (x_{1i} - \bar{x}_1) e_{1i} - \sum (x_{1i} - \bar{x}_1) (x_{2i} - \bar{x}_2) \sum (x_{2i} - \bar{x}_2) e_{1i}}{\sum (x_{1i} - \bar{x}_1) (x_{1i} - \bar{x}_1) \sum (x_{2i} - \bar{x}_2) (x_{2i} - \bar{x}_2) - [\sum (x_{1i} - \bar{x}_1) (x_{2i} - \bar{x}_2)]^2} \quad 2.42)$$

となる。仮定により $E(\sum (x_{1i} - \bar{x}_1) e_{1i}) = 0$ 、 $E(\sum (x_{2i} - \bar{x}_2) e_{1i}) = 0$ であるから、2.42) 式は

$$E(b_1) = b_1 \quad 2.43)$$

となるので、 b_1 は不偏推定量である。また一致推定量でもある。

ただし分散は 2.36) 式では

$V(\hat{b}_1) = \frac{\sigma^2}{\sum (x_{1i} - \bar{x}_1) (x_{1i} - \bar{x}_1)}$ であるが、2.37) 式では

$$V(\hat{b}_1) = \frac{\sigma^2}{(1 - r_{x_1 x_2}^2) \sum (x_{1i} - \bar{x}_1) (x_{1i} - \bar{x}_1)} \quad 2.44)$$

である。 $r_{x_1 x_2}$ は x_1 と x_2 の相関係数である。 $0 \leq r_{x_1 x_2} \leq 1$ であるから、 x_1 と x_2 が無相関でない限り、 $V(\hat{b}_1) > V(\hat{b}_1)$ となり、得られた分散の推定量は分散の最も小さい有効推定量ではない。このことは前に真のモデルが標準的線形回帰モデルの仮定を充たす場合、最小二乗法の推定量は最良線形不偏推定量 (BLUE) という結果からも示される。そのために過剰定式化したモデルを利用して仮説検定を行えば、 $H_0: b_1 = 0$ の帰無仮説を棄却しにくい (受容する) バイアスがかかることとなる。

(より一般的なモデルからより特定化されたモデルへ)

過小定式化の場合は得られた推計は一致性も不偏性もないので、その結果は全く用いることはできない。これに対し過剰定式化の場合は有効性はないが不偏性と一致性は保たれている。また不要な変数かどうかは t 検定や F 検定により、ある程度検証が可能である。それから考えれば、変数選択に迷ったときはその変数を入れる方が、致命的な誤りを避けるという意味でより望ましいといえ

よう。すなわちできるだけ制約のないモデルから出発し、不要な変数を順次除いていくことが望ましい。言い換えればより一般的なモデルから特定化されたモデルへ (general-to-specific) へ進むことが妥当である⁶⁾。

参考文献

仮説の検定の基礎的な考え方を紹介するものとしては、t分布やF文頭の各種分布の性質や導出を含めて

浅子・加納 [1998] 『入門経済統計学』(前出) の第5—8章

2.1 Kmenta. J [1986] *Elements of Econometrics: Michigan University*

のch 5 が分かりやすい。

ダミー変数や多重共線関係については

Maddala. G.S [1992] *Introduction To Econometrics* (前出) のCh 7—8 が詳しい。

6) f_{age} , $f_{age 2}$ を入れるケースはgeneralなケースと言えよう。この問題については後に時系列分析で更に触れることにする。