

## Logistic 回帰モデルにおける Wald 検定と TypeIII型検定

株式会社バイオスタティスティカル リサーチ

古川敏仁

Logistic 回帰モデルのパラメータ推定値  $\hat{\beta}$  は最尤法により推定される。最尤法では、対数尤度関数の漸近正規性より、パラメータ推定値の分散（標準誤差）が推定される。

今、 $p$  個の説明変数  $X=(x_1 x_2 \dots x_p)$  を考え、ロジット  $g(X)$  が下記のパラメータの線形式で示されるとする（Applied Logistic Regression second edition, Daivid W. Hosmer p31）。

$$g(X) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \quad (2.1)$$

このとき、 $\beta_j$  の推定値  $\hat{\beta}_j$  の統計学的有意性は、 $\beta_j$  の正規性を利用して、正規検定で簡単に行うことができる。

すなわち、 $W_j$  が 1.96 を超えれば、片側 5% で有意となる。これを Wald 検定と言う。

$$W_j = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)}$$

例えば、LOEBWT の RACE 変数をカテゴリとした、SAS 出力は以下のようになる。

```
PROC LOGISTIC DATA= LOEBWT;
```

```
CLASS RACE;
```

```
MODEL LOW= RACE /COBV;
```

```
RUN;
```

最大尤度推定値の分析						
パラメータ		自由度	推定値	標準誤差	Wald カイ 2 乗	Pr > ChiSq
Intercept		1	0.6612	0.1759	14.1255	0.0002
RACE	1	1	0.4936	0.2236	4.8716	0.0273
RACE	2	1	-0.3511	0.2889	1.4771	0.2242

RACE(1) に関しては、 $w=0.4936 \div 0.223=2.21$

SASでは、正規検定の代わりに、 $w_j$  を自乗した自由度 1 の  $\chi$  自乗検定の形で示されているため  $w^2=2.21^2=4.87$  となっている。

$$W_j^2 = \left( \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \right)^2$$

なぜ、正規検定の代わりに、 $\chi$  自乗検定の形で示しているかということ、それは、多変量への拡張性が良いからである。

今、個々のパラメータ推定値  $\beta_j$  の推定値  $\hat{\beta}_j$  の検定について示した。しかし、RACE は白人、黒人、その他の3つのカテゴリからなる変数であり、偶然 RACE(1)=黒人だけが、集団平均から離れている可能性がある (SAS の Logistic プロシジャでは、CLASS で指定された変数は deviation from means coding :Effect coding であり、その解釈には注意が必要である)。RACE 変数全体として、有意な情報を持っているかの検定も重要となる。これが Type III型の検定である。先ほどのプログラムの出力は、自由度 2 で  $p=0.0854$  で、RACE は有意とはなっていない。

Wald	Pr > ChiSq
カイ 2 乗	自由度
4.9209	2
	0.0854

この  $\chi$  自乗値はどのように計算したのであろうか。もし、 $p$ 個のパラメータが独立であれば、その  $W_j^2$ の総和は自由度  $p$ の  $\chi$  自乗分布に従う。

ためしに、RACE (1)と RACE (2)の  $\chi$  自乗値を足してみると、 $4.8716+1.4771 \approx 5.35 \neq 4.92$  となり、Type III型の  $\chi$  自乗値とは一致しない。これは、RACE (1)、RACE (2)が背反なカテゴリ変数であり集団平均からの差を求めるため、 $\hat{\beta}_1=0.4936$  と  $\hat{\beta}_2=-0.3511$  とは独立ではないことを示している。

Applied Logistic Regression second edition, Daivid W. Hosmer p39 では、 $p$  個の変数の独立成分の Wald  $\chi$  自乗統計量の総和 ( $p$  自由度) を求める式が書かれている。

$$W = \hat{\beta}' [\text{Var}(\hat{\beta})]^{-1} \hat{\beta} = \hat{\beta}' (X' V X)^{-1} \hat{\beta}$$

このような、行列表現は、最もシンプルな場合、つまり  $p=2$  で考えると理解しやすい。

今、 $\text{Var}(\hat{\beta})$ の逆行列を  $\phi_{11}$ 、 $\phi_{12}$ 、 $\phi_{21}$ 、 $\phi_{22}$  とおいてやって、実際の計算は何をやっているのかと確認すると

$$\begin{aligned} W &= \hat{\beta}' [\text{Var}(\hat{\beta})]^{-1} \hat{\beta} = \begin{bmatrix} \beta_1 & \beta_2 \end{bmatrix} \begin{bmatrix} \phi_{11} & \phi_{12} \\ \phi_{21} & \phi_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \beta_1 \phi_{11} + \beta_2 \phi_{21} & \beta_1 \phi_{12} + \beta_2 \phi_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \\ &= \beta_1 \beta_1 \phi_{11} + \beta_1 \beta_2 \phi_{21} + \beta_1 \beta_2 \phi_{12} + \beta_2 \beta_2 \phi_{22} \\ &\text{もし共分散がなければ} \\ &= \frac{\beta_1^2}{\text{Var}(\beta_1)} + \frac{\beta_2^2}{\text{Var}(\beta_2)} \end{aligned}$$

すなわち、もし、 $\beta_1$  と  $\beta_2$  が独立（共分散がなければ）であれば、個々のパラメータ推定値の Wald  $\chi^2$  自乗の和を求めることになる。ただ、現実はそのようなので、共分散分だけ調整した W 統計量を求めることになる。共分散が存在すれば以下を計算することになる。

$$W = \frac{\text{Var}(\beta_2)\beta_1^2}{(\text{Var}(\beta_1)\text{Var}(\beta_2) - \text{Cov}^2(\beta_1, \beta_2))} + \frac{\text{Var}(\beta_1)\beta_2^2}{(\text{Var}(\beta_1)\text{Var}(\beta_2) - \text{Cov}^2(\beta_1, \beta_2))} + \frac{-\text{Cov}(\beta_1, \beta_2)\beta_1\beta_2}{(\text{Var}(\beta_1)\text{Var}(\beta_2) - \text{Cov}^2(\beta_1, \beta_2))} \quad (1)$$

この手のことは、実際に計算してみると理解が早い。SAS によるパラメータ推定値の分散共分散行列は下記である。

推定共分散行列			
パラメータ	Intercept	RACE1	RACE2
Intercept	0.03095	-0.01189	0.021574
RACE1	-0.01189	0.050008	-0.04063
RACE2	0.021574	-0.04063	0.083475

RACE の Type III 型検定に必要な分散共分散行列は、その中の下記である。

パラメータ	RACE1	RACE2
RACE1	0.050008	-0.04063
RACE2	-0.04063	0.083475

また、必要なパラメータ推定値は下記である。

パラメータ	$\beta$ 推定値	
RACE	1	0.4936
RACE	2	-0.3511

これを、EXCEL で(1)に入れると  $W=4.92$  となる。

また、SAS には IML という行列計算言語があり、この手の計算がわかりやすく計算できる。下記のように、それぞれ行列に実際の値を入れ、V の逆行列を INV 関数で求めると、以下のような出力が得られる。

```
PROC ML;
RESET PRINT;
```

RESET IOG;

b={0.4936 -0.3511};

v={0.050008 -0.04063,-0.04063 0.083475};

gv=inv(v);

w=b\*gv\*b`;;

### SAS システム

B	1 row	2 cols	(numeric)
	0.4936		-0.3511

V	2 rows	2 cols	(numeric)
	0.050008		-0.04063
	-0.04063		0.083475

GV	2 rows	2 cols	(numeric)
	33.077472		16.099883
	16.099883		19.815972

W	1 row	1 col	(numeric)
			4.9214513

## まとめ

- Logistic 回帰モデルの Type III 型検定は Wald 統計量を用いている。
- Wald 統計量はパラメータ推定値が正規分布することを利用している。
- Wald  $\chi$  自乗値は、個々のパラメータに関する  $\chi$  自乗値を独立になるように調整し、その総和を求めたものである。ゆえに、 $p$  個の Wald  $\chi$  自乗値の自由度は  $p$  である。
- $W = \hat{\beta}' [\text{Var}(\hat{\beta})]^{-1} \hat{\beta} = \hat{\beta}' (X' V X) \hat{\beta}$  のような表現は一見難しそうに見えるが、行列は  $2 \times 2$ 、共分散なしで考えると、計算の意味が理解しやすい。
- また、今回の例のように、実際に計算すれば、案外、慣れるのは早い。この行列式が理解できるだけでも、SAS などのマニュアルのかなりの部分が理解可能となる。分からなくなったら何度も計算することが理解のためには重要である。