統計的方法への疑問

Rで計算して解決



2008.3 (生態学会大会、福岡)

統計的方法

『○●は、××のときに使うのはよくない』

『×△のときには、 ●○は 使うといい』

統計学の研究テーマでもある



誤ってはいないが抜けている

おおざっぱで役立たず

"誤った信仰"

目的変数が割合のとき直線回帰はまずい

ブロックは作れる限り作った方がいい

検定で有意にならなくても共変量としていれるといい

正規分布&等分散だったらt検定

分割表はカイ2乗検定

分割表は小さい値があればFisherの検定

割合のデータ解析ではoverdispersion無視はありえない

『分散に大きな差があるとき は、ノンパラメトリクス』

『ノンパラメトリクスを使えば 分布は気にしないでいい』

ほとんど都市伝説

解析的に解決すれば見通しもいいが・・・

解析的に解決しないことも多い

実際に計算して判断しよう

統計学の論文でも数値的に計算しているものが沢山



Rでやる"ご利益"

関数の計算結果の再利用

Rの豊富な関数たち やりたいことのかなりをカバー

すでにある関数 働かせる



R:組み合わせでできることが広がる

乱数

目次

統計的方法の代表的な"良い・悪い"

とても簡単な例. 相関係数rの検定

かなり簡単な例。ロジスティック回帰的な状況

モデル選択 (予測の最適化)

AIC 赤池情報量基準

-2×最大対数尤度+2×自由パラメーター数

正しいモデルの採用率

検定よい検定・悪い検定

第1種の誤り

有意確率・危険率に対応

本当は帰無仮説が正しいとき 帰無仮説を捨て対立仮説を採用

第2種の誤り

(1-検出力)

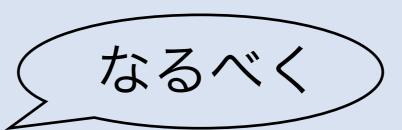
本当は対立仮説が正しいとき帰無仮説を捨てない

検定

よい検定・悪い検定

第1種の誤りの率が宣言した有意水準以下

のものの中で



第2種の誤りが小さい(検出力が高い)

Neyman-Pearson基準

正規分布のときの相関係数r

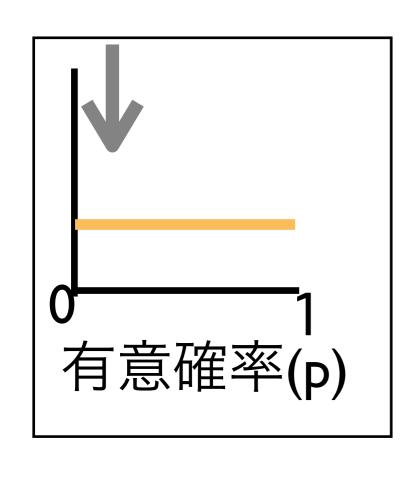
検定としての"良い・悪い" 第1段階

有意なものが出る率は宣言した水準以下か? 5%水準と言ってるのに10%水準とかに なっていないか

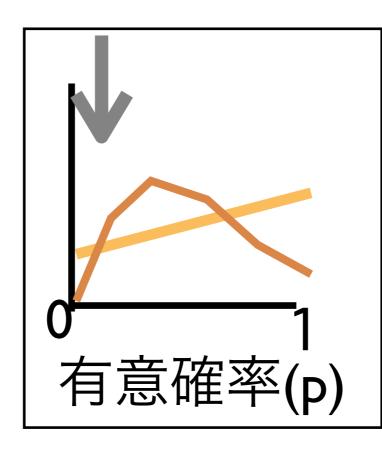
帰無仮説のモデルのもとで有意確率の分布

有意確率の分布

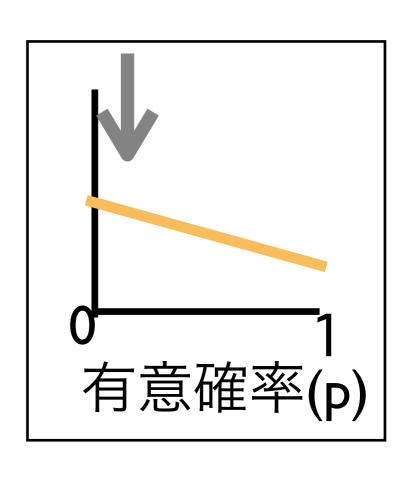
理想的には一様分布



理想的



鈍感



だめ (公称より甘い)

相関係数rの検定での有意確率の分布

(母集団:無相関、正規分布) サンプル数10

```
50000回繰り返し
pr1<-numeric(50000)
                     数値を入れておく場所
xx1<-numeric(10)
                       (ベクトル) 用意
yy1<-numeric(10)
for (i1 in 1:50000)
xx1<-rnorm(10,mean=9,sd=2)
yy1 < -rnorm(10, mean=7, sd=4)
rr1<-cor.test(xx1,yy1)
```

繰り返し計算

sum(pr1 <= 0.05)sum(pr1 <= 0.01)

pr1[i1]<-rr1\$p.value

計算結果集計

numeric(個数) 個数だけ数値の入れ場所を用意」

pr1<-numeric(50000)
xx1<-numeric(10)
yy1<-numeric(10)</pre>

10 サンプルサイズ (データ点個数)

50000 繰り返し数

数値入れ場所用意

繰り返し計算

計算結果集計

for (変数in範囲){} 変数を 1 ずつ変化させ繰り返し

```
for (i1 in 1:50000) {
    xx1<-rnorm(10, mean=9, sd=2)
    yy1<-rnorm(10, mean=7, sd=4)
    rr1<-cor.test(xx1, yy1)
    pr1[i1]<-rr1$p.value
}</pre>
```

```
xx1<-rnorm(10,mean=9,sd=2)
```

母平均9、母標準偏差2の正規乱数10個をxx1に入れる

```
for (i1 in 1:50000) {
    xx1<-rnorm(10, mean=9, sd=2)
    yy1<-rnorm(10, mean=7, sd=4)
    rr1<-cor.test(xx1, yy1)
    pr1[i1]<-rr1$p.value
}</pre>
```

```
rr1<-cor.test(xx1,yy1)
rr1に相関係数の検定結果
pr1[i1]<-rr1$p.value
rr1から有意確率を取り出してpr1に
```

```
for (i1 in 1:50000) {
    xx1<-rnorm(10, mean=9, sd=2)
    yy1<-rnorm(10, mean=7, sd=4)
    rr1<-cor.test(xx1, yy1)
    pr1[i1]<-rr1$p.value
}</pre>
```

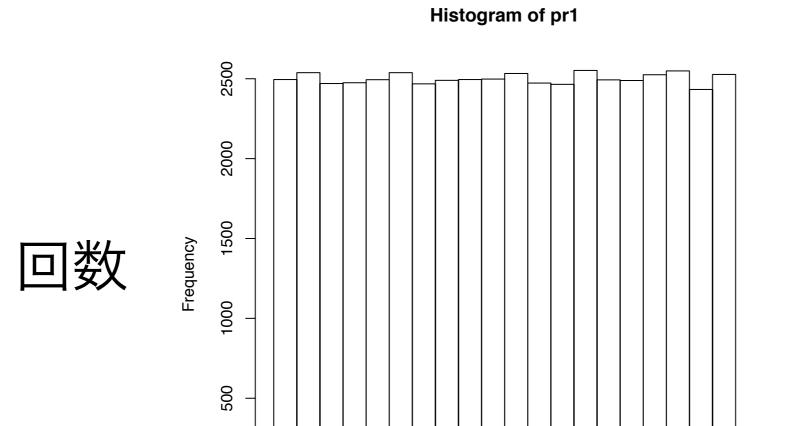
sum() 合計を計算する

sum(条件式) 条件式に合うデータの個数

$$sum(pr1 <= 0.05)$$

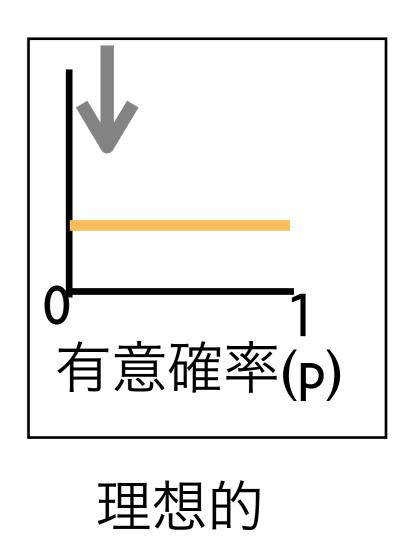
 $sum(pr1 <= 0.01)$

帰無仮説が正しい(母集団が無相関)とき



0.0

0.2



相関係数rを検定したときの有意確率

pr1

0.6

8.0

0.4

hist(pr1)

hist(pr1)で作ったものをPostscript保存してPDF化



どのくらいのばらつきが生じる?

5% 二項分布 分散

10000回 ±0.4%程度に0.95の確率で入る

1000回 ±1.4%程度に0.95の確率で入る

50% 二項分布 分散

10000回 ±1%程度に0.95の確率で入る

1000回 ±3%程度に0.95の確率で入る

ある条件

確率的なばらつき

計算上の表現

乱数

xx1<-rnorm(10,mean=9,sd=2)

乱数 → r+norm ← 正規分布

- d確率密度関数
- q 分布関数 P 確率点

正規分布でなく 対数正規分布 rnorm→rlnorm lnorm cauchy exp gamma f hbinom

出てきそうな確率分布は そろっている



相関係数rの検定での有意確率の分布

(母集団:無相関、正規分布) サンプル数10

```
50000回繰り返し
pr1<-numeric(50000)
                     数値を入れておく場所
xx1<-numeric(10)
                       (ベクトル) 用意
yy1<-numeric(10)
for (i1 in 1:50000)
xx1<-rnorm(10,mean=9,sd=2)
yy1 < -rnorm(10, mean=7, sd=4)
rr1<-cor.test(xx1,yy1)
pr1[i1]<-rr1$p.value
```

繰り返し計算

sum(pr1 <= 0.05)sum(pr1 <= 0.01)

計算結果集計

rr1<-cor.test(xx1,yy1)

```
> rr1
    Pearson's product-moment correlation
(中略)
t = 0.1516, df = 48, p-value = 0.8801
(中略)
sample estimates:
cor 0.02187758
```

rrlから、 ほしいもの(ここでは有意確率)を取り出す。



計算結果オブジェクトの中身調べ

関数attributes()

```
> attributes(rr1)

$names

[1] "statistic" "parameter" "p.value" "estimate"
"null.value"

[6] "alternative" "method" "data.name"
"conf.int"
```

```
> rr1$p.vaue [1] 0.880131
```



計算結果オブジェクトの中身調べ

関数str()

```
> str(rr1)
List of 9
$ statistic: Named num 0.152
 ..- attr(*, "names")= chr "t"
$ parameter : Named num 48
 ..- attr(*, "names")= chr "df"
$ p.value : num 0.88
$ estimate: Named num 0.0219
 ..- attr(*, "names")= chr "cor"
$ null.value : Named num 0
 ..- attr(*, "names")= chr "correlation"
$ alternative: chr "two.sided"
$ method : chr "Pearson's product-moment correlation"
$ data.name : chr "xx1 and yy1"
$ conf.int : atomic [1:2] -0.258 0.298
 ..- attr(*, "conf.level")= num 0.95
- attr(*, "class")= chr "htest"
```

rr1<-cor.test(xx1,yy1)

rrlから、有意確率を取り出す

> rr1\$p.vaue [1] 0.880131





計算結果オブジェクトの中身調べ

attributes()やstr()で調べて

計算結果オブジェクト名\$項目の名前

新しい関数では

計算結果オブジェクト名@項目の名前

S4オブジェクト

スロット

相関係数rの検定での有意確率の分布

(母集団:無相関、正規分布) サンプル数10

```
50000回繰り返し
pr1<-numeric(50000)
                     数値を入れておく場所
xx1<-numeric(10)
                       (ベクトル) 用意
yy1<-numeric(10)
for (i1 in 1:50000)
xx1<-rnorm(10,mean=9,sd=2)
yy1 < -rnorm(10, mean=7, sd=4)
```

繰り返し計算

sum(pr1 <= 0.05)sum(pr1 <= 0.01)

rr1<-cor.test(xx1,yy1)

pr1[i1]<-rr1\$p.value

計算結果集計

相関係数rの検定

乱数は使っているが 計算できない場合が生じない

> xが全部同じ値 yが全部同じ値

乱数+繰り返し計算 予想しない数値 エラーが出て途中で止まる

関数try() 関数tryCatch()



目次

統計的方法の代表的な"良い・悪い"

とても簡単な例. 相関係数rの検定

かなり簡単な例。ロジスティック回帰的な状況



体サイズ 温度 標高 など

生存/死亡 メス/オス 勝ち/負け など

ロジスティック回帰 量的変数が確率に影響

2 処理の比較 生存個体と死亡個体の体サイズ 原因/結果が逆転?



体サイズ 温度 標高 など



2値変数

生存/死亡 勝ち/負け メス/オス など

ロジスティック回帰

尤度比検定

Wald検定

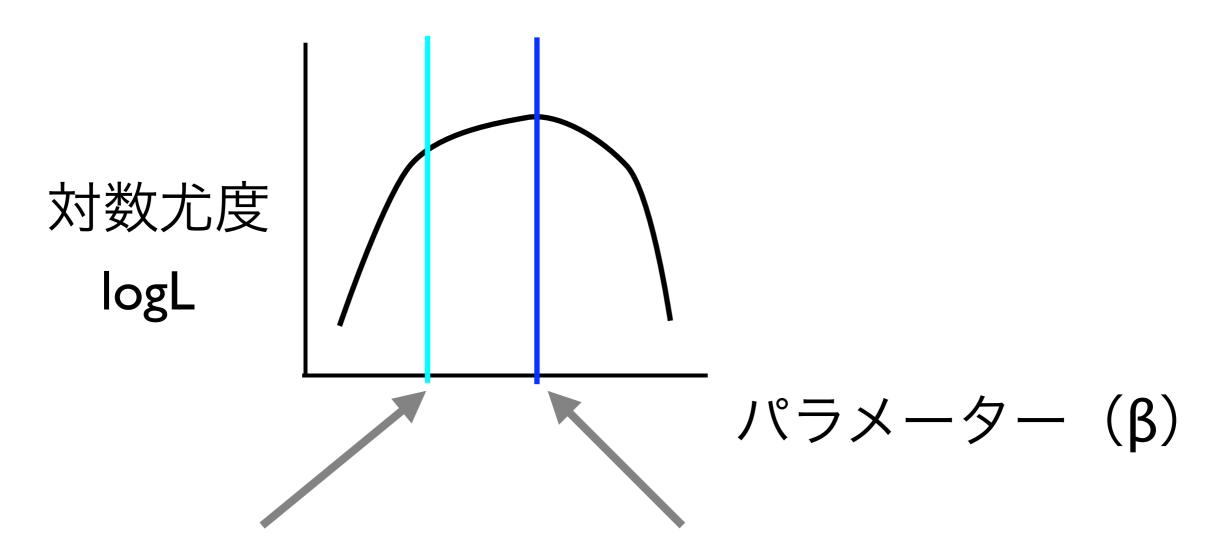
2処理の比較

中央値検定

Wilcoxon検定 (Mann-WhitneyのU検定)

t検定(Welch)

最尤法(maximum likelihood)



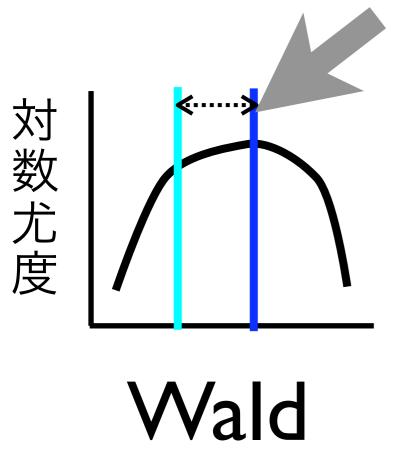
帰無仮説のモデル

最尤推定值

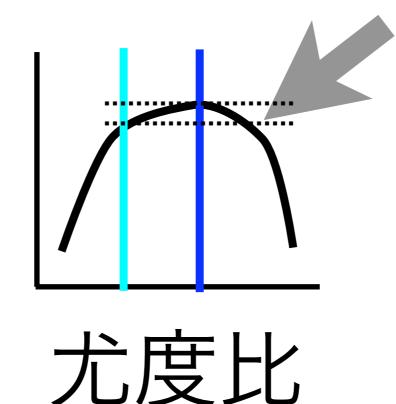
直線回帰なら0

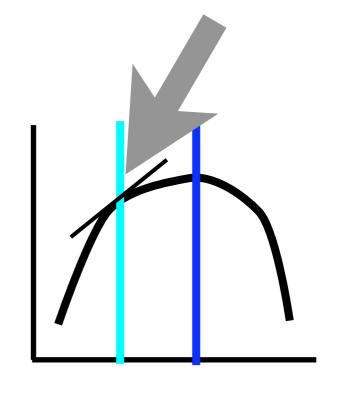
直線回帰なら 推定された傾き

最尤法と3つの検定



最尤推定值





対数尤度の差

スコア

(Lagrange乗数)帰無仮説のモデルでの対数尤度の傾き

サンプル数無限大では同じ

量的変数 体サイズ 温度 標高 など



2値変数 生存/死亡 メス/オス など

ベルヌイ過程(overdispersionなし)

サンプルサイズ100 (1000までやった)

使った プログラム

やや長いの

はいろいろ

な方法を評

価したから

```
fpI < -numeric(10000)
wp I <-numeric(10000)
lpI < -numeric(10000)
tp I <-numeric(10000)
Irpl<-numeric(10000)
cntl<-numeric(10000)
fcntl<-numeric(10000)
aic I <-numeric (10000)
n1<-100
xI < -numeric(nI)
for (i in 1:n1) {
xI[i]<-qnorm(i/(nI+I),mean=20.0,sd=4.0)
dI<-numeric(nI)
for (i in I:n1) {
dI[i] < -0.5
yl<-numeric(nl)
for (il in 1:10000) {
for (i in 1:n1) {
if (d1[i]>=runif(1)) {y1[i]<-1} else {y1[i]<-0}
xa < -sum(head(y | 1, n | 1/2))
xc<-(n1/2-xa)
xb < -sum(tail(y1,n1/2))
xd < -(n1/2-xb)
fm I < -matrix(c(xa,xb,xc,xd),ncol=2)
fp | [i | ] <- fisher.test(fm | ) $p. value
fcntl[il]<-(xa+xb+xc+xd)
yfl<-factor(yl)
wwl < -wilcox.test(xl \sim yl)
wpl[il]<-wwl$p.value
lgt0 < -glm(yfl \sim I, family = binomial)
|gt| < -glm(yfl \sim x l, family = binomial)
slgt I <- summary(lgt I)
lp | [i | ] <-slgt | $coef[2,4]</pre>
|\text{Irt}| < 2*abs(|\text{logLik}(|\text{gt}|)[1]-|\text{logLik}(|\text{gt}0)[1])
lrpl[il] < -(l-pchisq(lrtl,l))
trl < -t.test(xl \sim yfl)
tpl[il]<-trl$p.value
cntl[il]<-l
aic [i1]<-(lgt | $aic-lgt0$aic)
sum(f_D I <= 0.05)
sum(lp I <= 0.05)
sum(wp I <= 0.05)
sum(Irp I <= 0.05)
sum(tp I <= 0.05)
sum(fpl <= 0.01)
sum(lp1 <= 0.01)
sum(wpl <= 0.01)
sum(IrpI <= 0.01)
sum(tpl <= 0.01)
sum(aicl <= 0.0)
sum(fcnt1)/10000
```

尤度比、

Wald

```
lp1 < -numeric(10000)
               lrp1<-numeric(10000)
               aic1<-numeric(10000)
               n1<-100
               x1<-numeric(n1)
AICのみ for (i in 1:n1) {
               x1[i] < -qnorm(i/(n1+1), mean = 20.0
                                                     sd=4.0
               d1<-numeric(n1)
               for (i in 1:n1) {
               d1[i] < -0.5
               yl<-numeric(n1)
               for (i1 in 1:10000) {
               for (i in 1:n1) {
               if (d1[i] > = runif(1)) \{y1[i] < -1\} else \{y1\}
                                                         返し
               yfl<-factor(y1)
               lgt0<-glm(yf1~1,family=binomial)
               lgt1<-glm(yf1~x1,family=binomial)
               slgt1<-summary(lgt1)
               lp1[i1]<-slgt1$coef ["x1","Pr(>|z|)"]
lrt1<-2*abs(logLik(lgt1)[1]-logLik(lgt0
                                                        [1]
               lrp1[i1]<-(1-pchisq(lrt1,1))
               aic1[i1]<-(lgt1$aic-lgt0$aic)
               sum(lp1 <= 0.05)
                                            計算結果集計
               sum(lrp1 <= 0.05)
               sum(aic1 \le 0.0)
```



glm()の結果オブジェクトから検定結果を取り出す (例、有意確率)

attributes()やstr()でさがす

> glmの結果のsummaryオブジェクト\$coefficients

Estimate Std. Error z value Pr(>|z|) (Intercept) 18.23579822 3.967524 4.5962662 3.140334e-05 x01 0.06201217 0.409028 0.1516086 8.801311e-01

glmの結果のsummaryオブジェクトのcoefficients 行や列に名前のついた行列

coefficients["x01","Pr(>|z|)"]

↑久保さん指摘による

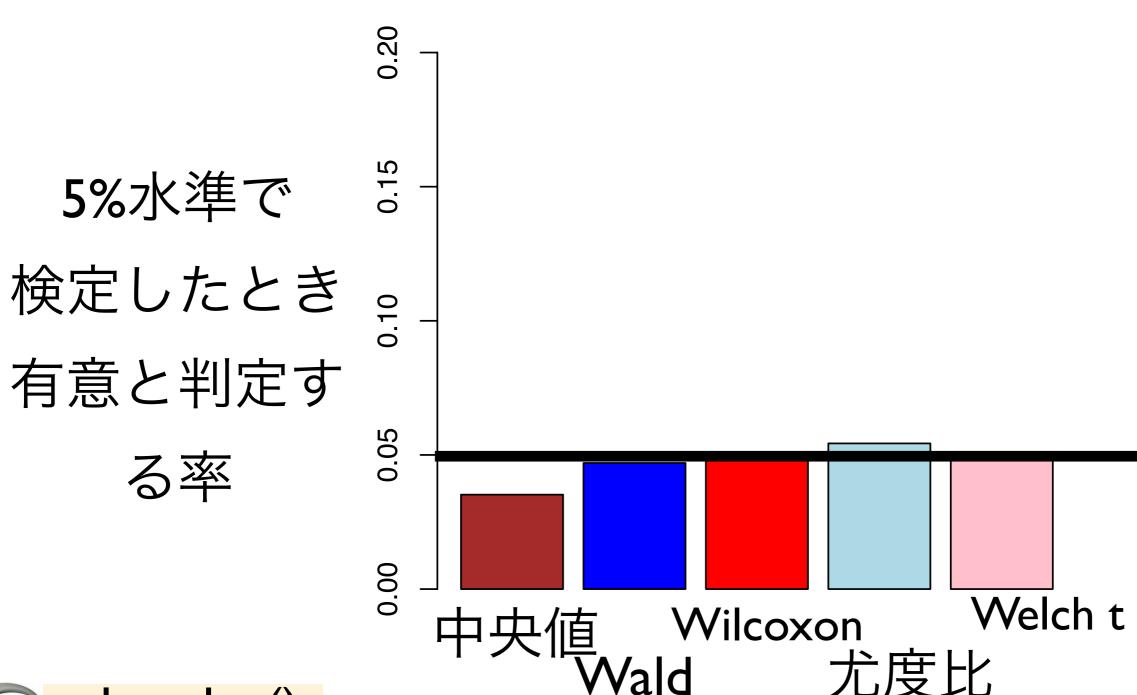
あるいは

glmの結果のsummaryオブジェクトcoefficients[8] または

glmの結果のsummaryオブジェクト\$coefficients[2,4]

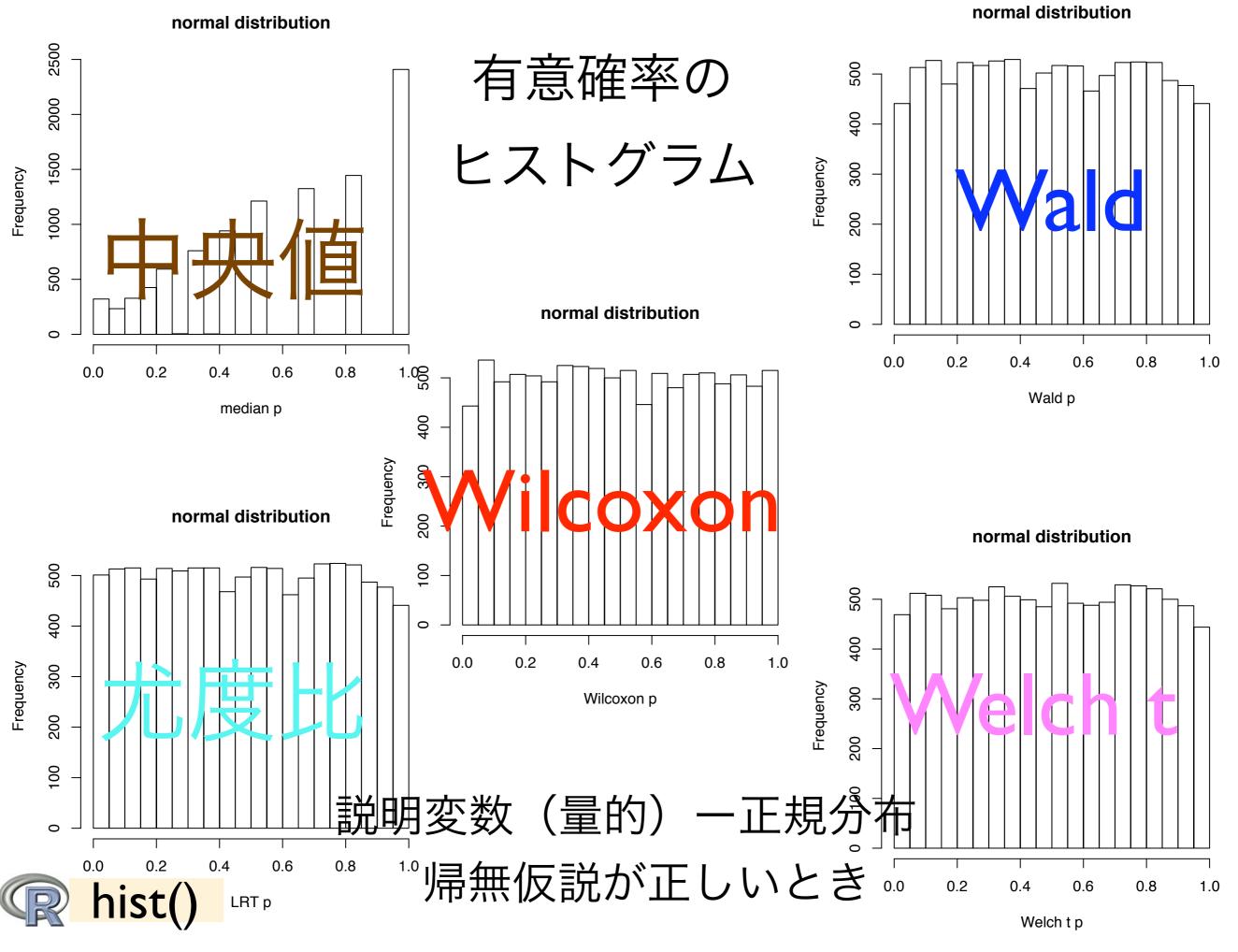
説明変数 (量的) 一正規分布 帰無仮説が正しいとき

normal distribution



barplot()

type 1 error rate



説明変数の分布も解析結果に影響(目的変数の分布は注意を払われるが)

説明変数一目的変数をひっくり返すのは要注意

検定specific:単に第1種の誤りの 率が宣言値以下かだけでは不十分

正規線形モデル(誤差=等分散の正規分布)では

エラーが出て途中で止まる

計算できない場合 目的変数が全部同じ値

x > ある値 y がすべて I 完全分離 x < ある値 y がすべて 0

関数try() 関数tryCatch()

関数の計算結果の再利用

その他の場合にも

面倒なところはRの既存の関数で

例、パラメトリック・ブートストラップ