

2018年度、統計関連学会連合大会  
9月9日～13日

# 単変量Cox回帰にもとづく遺伝子選択と 複合共変量による生存期間の予測

発表者；江村剛志（国立中央大学，統計研究所）  
松井茂之（名古屋大学大学院、医学部）  
Hsuan-yu Chen (中央研究院、統計科学研究所)

本発表の詳細な解説は大会ウェブページに掲載される「詳細論文」  
または次の原論文

**Emura T, Matsui S, Chen HY (2019)**, compound.Cox: univariate feature selection and compound covariate for predicting survival, in revision *Computer Methods and Programs in Biomedicine*.

# 研究の設定 & 動機

{ 生存期間 = 応答変数  
遺伝子発現量 = 共変量 (予後因子)

大量の遺伝子の中から、生存期間に関連する遺伝子を選択したい

- 肺がん患者  
*ERBB3, LCK, DUSP6, STAT2* (Chen et al., 2006 NEJM)
- 乳がん患者  
*ECRG4* (Sabatier et al., 2011, PLoS ONE)
- 卵巣がん患者  
*CXCL12* (Pople et al., 2012, British J. of Cancer)

# データの解析例

- 肺がん患者

*ERBB3, LCK, DUSP6, STAT2*

など16の遺伝子が生存期間に関連

(Chen et al., 2007 NEJM)

☞ データから単変量Cox回帰に基づいて、  
16の関連遺伝子を選択。

使用したデータ;

$n=125$ 人の肺がん患者(台湾の大学病院)  
(死亡38人 + 打ち切りによる生存87人)

$p=485$ 個の遺伝子

(選択された関連遺伝子は16個)

# 肺がんデータ in

## 『*compound.Cox*』 Rパッケージ (Emura et al. 2018)

訓練標本 ( $n=63$ ) または  
テスト標本 ( $n=62$ ) の指標

```
> library(compound.Cox)
> data("Lung") 打ち切り
> Lung
```

```
→→ t.vec      d.vec  train  VHL  IHPK1  ...  RPL5
1    47.06271    0    FALSE  2     2           4
2    49.27393    0     TRUE  3     4           4
3    20.06601    1     TRUE  2     3           1
4    26.99670    1     TRUE  2     4           2
5    39.90099    0    FALSE  3     4           4
⋮      ⋮      ⋮      ⋮      ⋮      ⋮      ⋮
125  56.84141    0    FALSE  3     2     ...    3
```

↓  $p=97$ 個の遺伝子 (離散値)

Chen et al. (2007) では

訓練標本 ( $n=63$ ) に単変量Cox回帰を用いて関連遺伝子を選択、  
テスト標本 ( $n=62$ ) で関連遺伝子の予測能力を評価

$T$  = 死亡時間

$x_j$  =  $j$ 番目の遺伝子発現量

$T$  と  $x_j$  関連強  $\rightarrow$  単変量Cox回帰の推定値が大

$$h(t | x_j) = \frac{\Pr(t \leq T < t + dt | T \geq t, x_j)}{dt} = h_0(t) \exp(\beta_j x_j)$$

## データの形式

$\{ (t_i, \delta_i, \mathbf{x}_i); i = 1, \dots, n \}$ ,

- $t_i$  : survival time or censoring time,
- $\delta_i$  : censoring indicator (  $\delta_i = 1$  if  $t_i$  is survival time, or  $\delta_i = 0$  if  $t_i$  is censoring time ),
- $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$ :  $p$ -dimensional features (genes).

# 単変量Cox回帰にもとづく遺伝子選択法

Step1: 単変量Coxモデルで全遺伝子を個々に検定(多重検定)

$$H_{0j} : \beta_j = 0 \quad \text{vs.} \quad H_{1j} : \beta_j \neq 0$$

$$h_{0j}(t) \exp(\beta_j x_{ij}), \quad j = 1, \dots, p$$

(1) ワルド検定 ; z値 =  $\hat{\beta}_j / SE(\hat{\beta}_j)$

$\hat{\beta}_j$  = 単変量部分尤度推定値

(2) スコア検定 ; z値 =  $S_j / V_j^{1/2}$  = 単変量スコア統計量 ÷ SD

Step2 : 棄却された遺伝子を選択

例 ; P値が0.05以下である遺伝子を選択

Step3 : 選択された遺伝子の

(1) FDR値を計算(簡便法、Permutation法)

(2) CVL値を計算(部分尤度のK分割クロスバリデーション)

Step4 : 選択された複数の遺伝子で生存期間の予後予測(後述)

本研究ではSteps 1-4を実行するRパッケージ、`compound.Cox` を提案

# P値0.05でRパッケージを訓練標本に適用、 遺伝子選択

```
> res=uni.Wald(t.vec,d.vec,X.mat)
```

```
> res$beta[res$P<0.05]
```

```
HMMR      LCK      ANXA5      IRF4      STAT2      ERBB3      NF1
0.5156711 -0.8447389 -1.0876762 0.5176704 0.5849869 0.5509026 0.4715235
DLG2      HGF      CPEB4      ZNF264     MMD      RNF4      FRAP1
1.3215044 0.5086750 0.5891676 0.5473276 0.9151541 0.6463635 -0.7696768
STAT1      DUSP6
-0.5844262 0.7524497
```

↑ [Chen et al. \(2007 NEJM\)](#)と同じ結果

- この16の遺伝子の中に、偽陽性はいくつあるか？  
⇒FDR (False discovery rate) を計算

# False Discovery Rate (FDR)とは

FDR=選択された遺伝子中の無関係な遺伝子の割合  
(通常、20%以下であることが望ましい(?))

	選択された 遺伝子	選択されな かった遺伝子	全遺伝子
関連する遺伝子			
無関係な遺伝子	$f$		
	$q=16$		$p=97$

$$\text{FDR} = f/16$$

$$\text{FDR} = 0.05 \times 97/16 = 0.30 \quad (30\%とやや大きい)$$



# FDRに関する一般的注釈

- **Permutation法でも計算可** (Witten & Tibshirani 2010 *SMMR*)
  - *compound.Cox* パッケージで計算を実装
- **FDRは期待値;**  
実際の偽陽性の数や、どれが偽陽性なの不明
- **FDRは選択された遺伝子の予測能力ではない**  
(後述のCVL値を予測能力として使用)

# CVL (Cross-validated likelihood) とは

P値の閾値を与えたもとで選択された複数遺伝子の予測能力の指標で、部分尤度のK分割クロスバリデーションを用いて、次の式で定義

$$CVL = \sum_{k=1}^K \{ \ell(\hat{\gamma}_{-k}) - \ell_{-k}(\hat{\gamma}_{-k}) \},$$

where  $\hat{\gamma}_{-k} = \arg \max_{\gamma} \ell_{-k}(\gamma)$ ,

$$\ell(\gamma) = \sum_i \delta_i \left[ \gamma \text{CC}_{i,-k} - \log \left\{ \sum_{\ell \in R_i} \exp(\gamma \text{CC}_{\ell,-k}) \right\} \right],$$

$$\text{CC}_{i,-k} = \sum_{j \in \Omega_{-k}} w_{j,-k} x_{ij}$$

$$\ell_{-k}(\gamma) = \sum_{i \in \mathfrak{I}_{-k}} \delta_i \left[ \gamma \text{CC}_{i,-k} - \log \left\{ \sum_{\ell \in R_i \cap \mathfrak{I}_{-k}} \exp(\gamma \text{CC}_{\ell,-k}) \right\} \right],$$

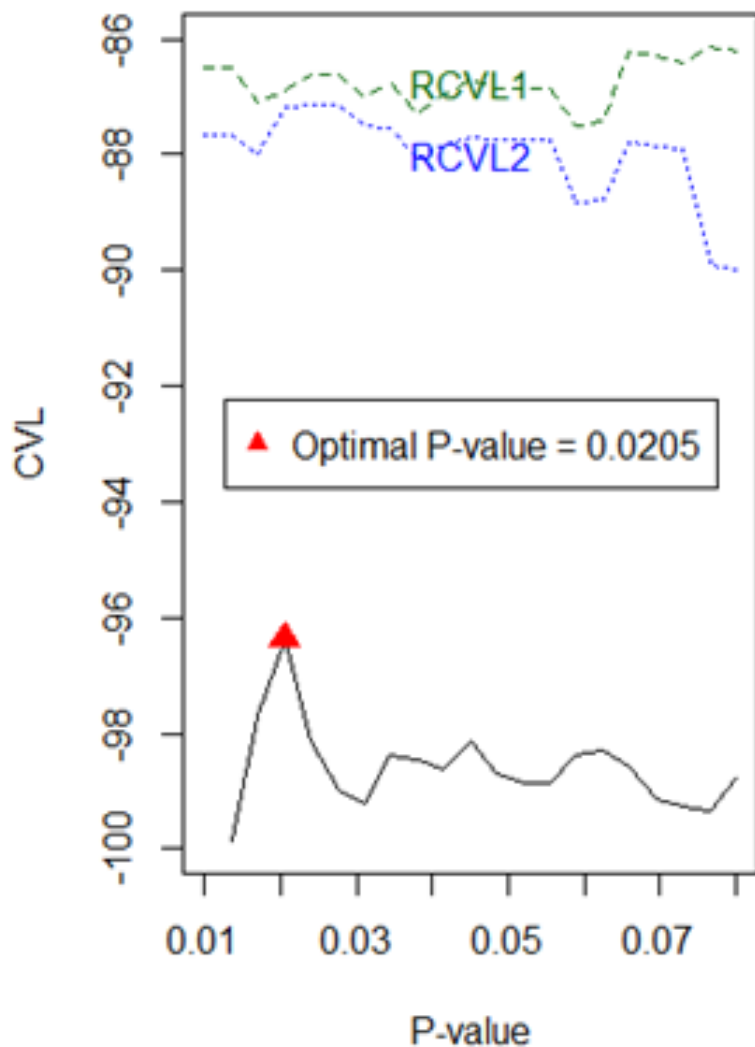
高CVL値  $\Rightarrow$  高予測能力; 詳細は下記文献

Matsui 2006; *BMC Bioinformatics*

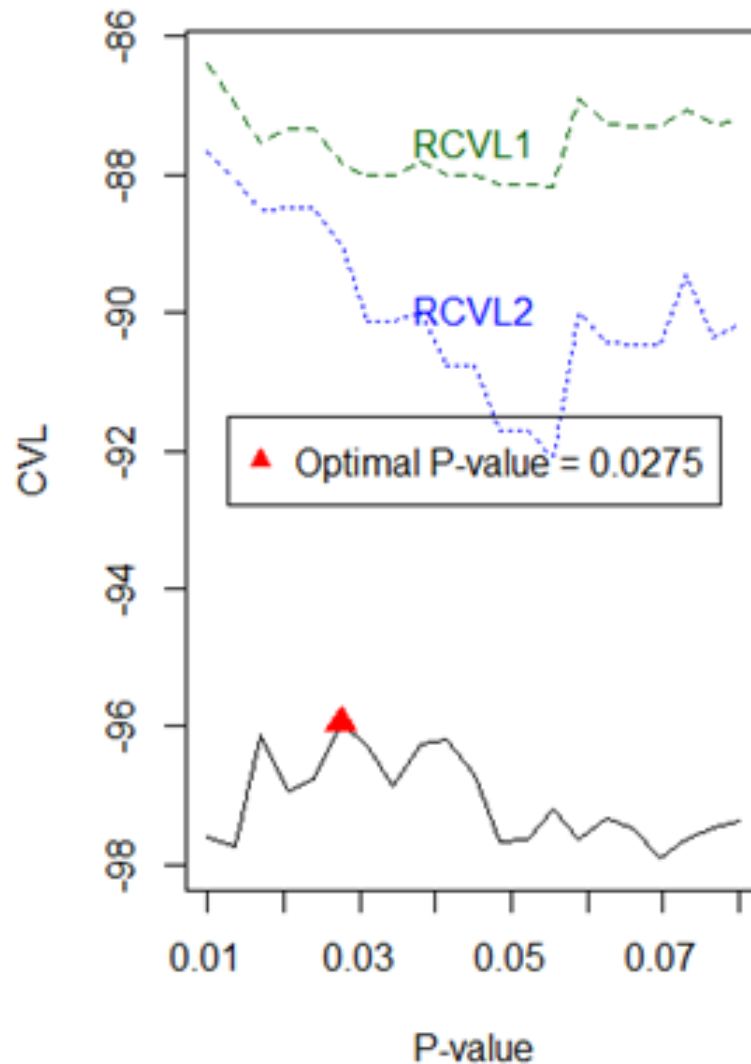
Emura, Matsui and Chen 2019 *Computer Methods and Programs in Biomed*

# CVL値を最適化するP値の閾値を求める

Wald test



Score test



# ワールド検定の最適閾値 (P = 0.0205) で 7個の遺伝子を選択

```
> uni.selection(t.vec,d.vec,X.mat,K=20, P.value=0.0205,score=FALSE)
```

```
$beta
```

ANXA5	DLG2	ZNF264	DUSP6	CPEB4	LCK	STAT1
-1.0876762	1.3215044	0.5473276	0.7524497	0.5891676	-0.8447389	-0.5844262

```
$CVL -96.37303
```

↑CVL値

FDR値=0.0205 × 97/7=0.29 (29%)

---

# スコア検定の最適閾値 (P = 0.0275) で 10個の遺伝子を選択

```
> uni.selection(t.vec,d.vec,X.mat,K=20, P.value=0.0275,score=TRUE)
```

```
$Z
```

ANXA5	DLG2	ZNF264	DUSP6	CPEB4	LCK	STAT1	STAT2
-3.363578	3.111772	2.814363	2.710854	2.538888	-2.511423	-2.445038	2.369334
RNF4	IRF4						
2.345912	2.231286						

```
⋮
```

```
$CVL -95.95690
```

↑CVL値

FDR値=0.0275 × 97/10=0.30 (30%)

# 生存期間の予後予測

- 選択された複数遺伝子の発現量;  $(x_1, \dots, x_q)$

スコア検定の最適閾値では  $q=10$

- 複合共変量 (Compound Covariate):

$$CC = w_1 x_1 + \dots + w_p x_q$$

$\beta$ 値;  $(w_1, \dots, w_q) = (\hat{\beta}_1, \dots, \hat{\beta}_q)$

$z$ 値;  $(w_1, \dots, w_q) = (z_1, \dots, z_q)$

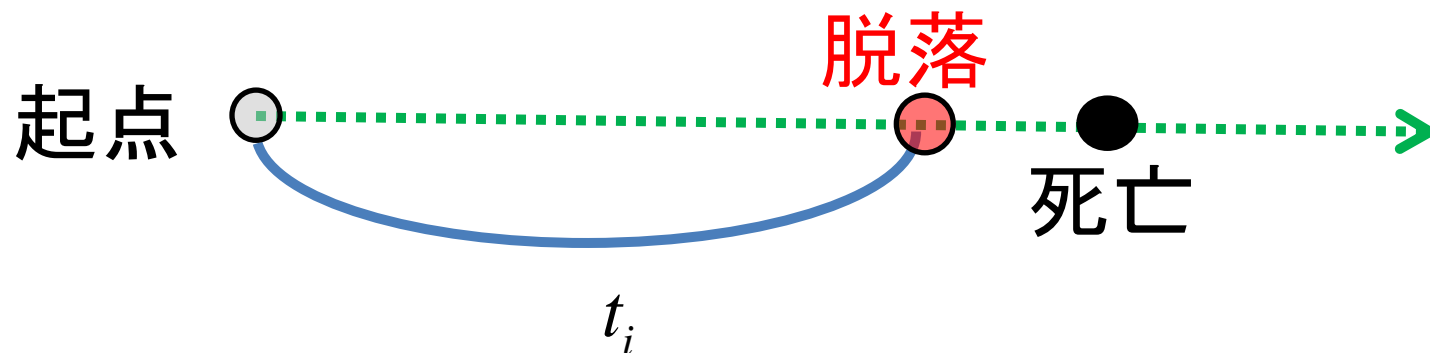
- 予後分類:  $CC < c \Rightarrow$  予後は良い  
 $CC > c \Rightarrow$  予後は悪い

ここで、 $c = CC$ の閾値

# 従属センサー (Dependent censoring) 問題

死亡の直前に**脱落**する患者

⇒センサリングと死亡に相関あり



独立センサーの仮定が崩れ、  
Cox回帰の推定値にバイアスが生じる  
(Emura & Chen 2016 *SMMR*)

# 従属センサーのコピュラ・モデル

$T$  = 生存期間(死亡時間)

$U$  = センサー(脱落)時間

$x_j$  =  $j$ 番目の遺伝子発現量

$C_\alpha$  = コピュラ関数  
( $\alpha$ は相関パラメータ)

↓ 2次元生存関数

$$\Pr(T_i > t, U_i > u | x_{ij}) = C_\alpha \{ \Pr(T_i > t | x_{ij}), \Pr(U_i > u | x_{ij}) \}$$

$$\Pr(T_i > t | x_{ij}) = \exp \{ -\Lambda_{0j}(t) e^{\beta_j x_{ij}} \}$$

遺伝子  $j$  が  $T$  に与える影響

SPRINGER BRIEFS IN STATISTICS  
JSS RESEARCH SERIES IN STATISTICS

Takeshi Emura  
Yi-Hau Chen

Analysis of Survival  
Data with Dependent  
Censoring  
Copula-Based  
Approaches



Springer

# コンピュータによる従属センサーへの対応

## セミパラメトリック最尤推定量 (Chen 2010, JRSSB)

$$\begin{aligned} & \ell(\beta_j, \gamma_j, \Lambda_{0j}, \Gamma_{0j} | \alpha) \\ &= \sum_i \delta_i [ \beta_j x_{ij} + \log \eta_{1ij}(t_i; \beta_j, \gamma_j, \Lambda_{0j}, \Gamma_{0j} | \alpha) + \log d\Lambda_{0j}(t_i) ] \\ &+ \sum_i (1 - \delta_i) [ \gamma_j x_{ij} + \log \eta_{2ij}(t_i; \beta_j, \gamma_j, \Lambda_{0j}, \Gamma_{0j} | \alpha) + \log d\Gamma_{0j}(t_i) ] \\ &- \sum_i \Phi_\alpha [ \exp\{ -\Lambda_{0j}(t_i) e^{\beta_j x_{ij}} \}, \exp\{ -\Gamma_{0j}(t_i) e^{\gamma_j x_{ij}} \} ], \end{aligned}$$

compound.Cox パッケージで計算

$$(\hat{\beta}_j(\alpha), \hat{\gamma}_j(\alpha), \hat{\Lambda}_{0j}(\alpha), \hat{\Gamma}_{0j}(\alpha))$$

予測に使う  $j$  番目の遺伝子の重み  $w_j$  を計算



# 複数遺伝子を使った生存予後予測

1. Optimal Wald (7個の遺伝子):

$$CC = \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_7 x_{i7}$$

2. Optimal score (10個の遺伝子):

$$CC = z_1 x_{i1} + \cdots + z_{10} x_{i10}$$

3. Optimal Wald + copula (7個の遺伝子):

$$CC = \hat{\beta}_1(\hat{\alpha}) x_{i1} + \cdots + \hat{\beta}_7(\hat{\alpha}) x_{i7}$$

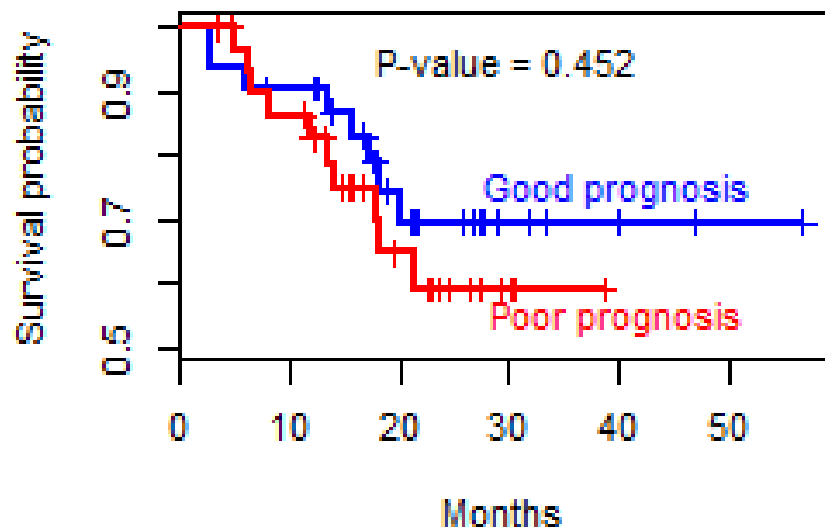
4. Optimal score + copula (10個の遺伝子)

$$CC = \hat{\beta}_1(\hat{\alpha}) x_{i1} + \cdots + \hat{\beta}_{10}(\hat{\alpha}) x_{i10}$$

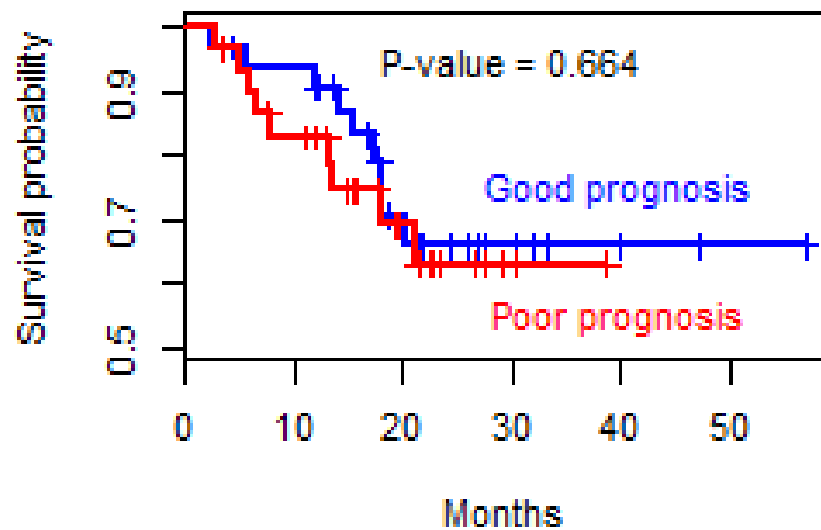
$CC < c \Rightarrow$  予後は良い(高生存率)

$CC > c \Rightarrow$  予後は悪い(低生存率)

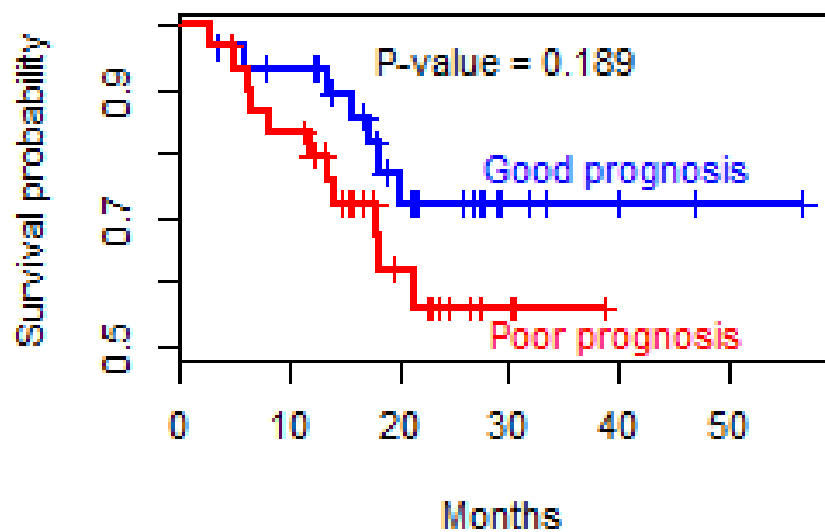
### Optimal Wald test



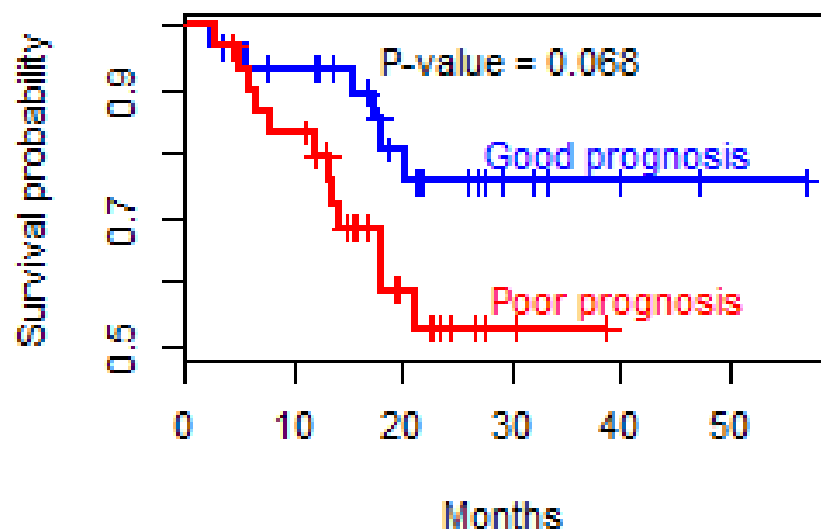
### Optimal score test



### Copula + optimal Wald test



### Copula + optimal score test



# 本研究のまとめ

- **Rパッケージ; `compound.Cox` を提案**
  - 単変量Cox回帰法にもとづく遺伝子選択を多重検定で実行 (Lassoなど罰則付尤度を使う流行りの手法と異なる)
- **選択された遺伝子の評価法を実装**
  - 予測能力評価 (CVL値) [Matsui \(2006 BMC Bioinformatics\)](#)
  - 無関係な遺伝子の割合 (FDR値)  
[Witten and Tibshirani \(2010 SMMR\)](#)
- **従属センサーの問題への対処法を実装**
  - 回帰推定値のバイアス補正 [Emura and Chen \(2016 SMMR\)](#)
  - コピュラ・グラフィック推定量による、生存関数のバイアス補正  
[Rivest and Wells \(2001 JMVA\)](#)

# 今後の課題

- **単変量Cox回帰の理論的裏付け**

なぜ単変量回帰が、多変量回帰よりも良いのか？

シミュレーションや実データで実証されてきたが。。。

- **死亡時間＋中間イベント時間が観測されるケース**

*compound.Cox* は利用可能であるが、中間イベントが死亡によって従属センサーされるなど問題点もある。

- **動的予測への応用；**（遺伝子）＋（中間イベント情報）

の同時利用で生存期間の予測精度を上げる

例；増悪あり・なしの影響を予測モデルに組み入れる。

Emura et al. (2018 *SMMR*) 医者や臨床試験に組み込むには。。。

- **平均余命の予測；**

パラメトリックモデルで患者の生存関数形を特定しているときの遺伝子選択は。。。

Wu, Michimae and Emura (2019- submitted)

Shinohara, Emura, Michimae, Takeuchi (2019 in preparation)