

# Introduction to Survival Analysis

Statistical Physics Approaches to Systems Biology, Havana, Feb 2019

ACC Coolen

King's College London and Saddle Point Science



## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

1982:  
Commodore 64



next generation data  
previous generation  
analysis ...

1982:  
Commodore 64



## Regression Models and Life-Tables

D. R. Cox

*Journal of the Royal Statistical Society. Series B (Methodological)*, Volume 34, Issue 2  
(1972), 187-220.

Stable URL:

<http://links.jstor.org/sici?sici=0035-9246%281972%2934%3A2%3C187%3ARMAL%3E2.0.CO%3B2-6>



## Introduction

### Post-genome medicine

Statistics in medicine is tricky

## The formalism of survival analysis

Terminology and objectives

Survival probability and cause-specific hazard rates

Data likelihood in terms of cause specific hazard rates

Pitfalls and misconceptions

Special cases

## Event time correlations and identifiability

Independently distributed event times

The identifiability problem

## Proportional hazards regression

Definitions and assumptions

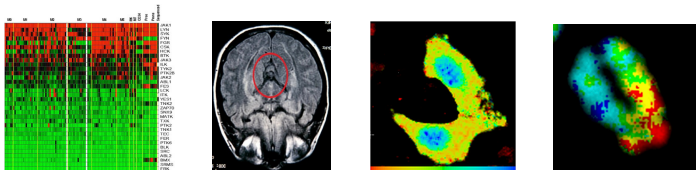
Parameter estimation from data

ML parameters of Cox's model

Uniqueness, error bars, and  $p$ -values

# Biomedicine has changed drastically in recent decades

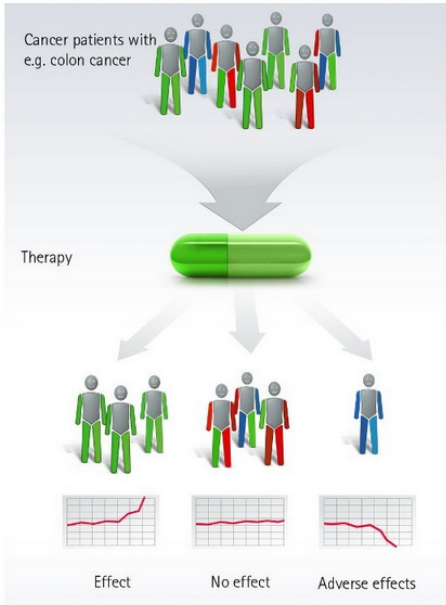
## modern biomedical data



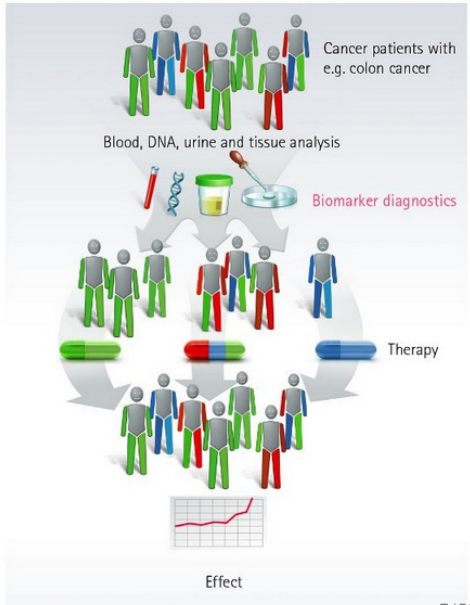
- ▶ *volume* of data ...
- ▶ *diversity* of data sources  
(clinical, genomic, biomarkers, health records, imaging, ...)
- ▶ *complexity* of experimental pipelines  
(confounders, batch effects, variability between centres, ...)
- ▶ *dimensionality* of measurements  
clinical images ( $\sim 10^6$ ), transcriptome/proteome ( $\sim 10^5$ ),  
DNA and methylation ( $\sim 10^{10}$ ), ...

# Personalized medicine: tailored treatments

Medicine of the present: one treatment fits all



Medicine of the future: more personalized diagnostics



generating 'big data' is not enough ...



- ▶ 'right drug, right dose, at the right time ...'

need *predictive models*  $p(y|\mathbf{z})$ ,

$\mathbf{z}$ : individual's makeup (DNA, gene expr, metabolism, environment, ...)

$y$ : response to treatment

- ▶ regression:

find parameters  $\theta$  of model  $p(y|\mathbf{z}, \theta)$  from historic data

curse of dimensionality ...

pre-genome medicine:  $N \sim 10^3$  data points,  $\dim(\theta) \sim 10^2$

post-genome medicine:  $N \sim 10^4$  data points,  $\dim(\theta) \sim 10^{10}$

- ▶ simpler question: predict patient's *individual risk*  
(target aggressive treatments wisely)

cancer, heart disease, diabetes, ...:

relevant outcome is often a *duration*  $t$ ,

OS (overall survival), or PFS (progression-free survival)

predictive model:  $p(t|\mathbf{z}, \theta)$



## Introduction

Post-genome medicine

Statistics in medicine is tricky

## The formalism of survival analysis

Terminology and objectives

Survival probability and cause-specific hazard rates

Data likelihood in terms of cause specific hazard rates

Pitfalls and misconceptions

Special cases

## Event time correlations and identifiability

Independently distributed event times

The identifiability problem

## Proportional hazards regression

Definitions and assumptions

Parameter estimation from data

ML parameters of Cox's model

Uniqueness, error bars, and  $p$ -values

Statistics in medicine:  
tricky business ...



"I can prove it or disprove it! What do you want me to do?"

why is statistics tricky?  
Monty Hall problem

'Let's Make a Deal'  
(USA gameshow, 1963-1977)



standard quiz show,  
winner has to choose prize at the end,  
three doors: one with big prize, two with goats ...



- winner selects one door randomly
- Monty opens one door with a goat (two doors left ...)
- Monty gives winner the chance to change selection at last minute

*would it matter?*

## The main pitfalls in statistics

- ▶ *accidental conditioning*

(Monty Hall problem, share statistics,  
shop opening hours consultation, ...)

extra knowledge:

→ *reduces possibilities*

→ *affects probabilities*

$$\overbrace{P(A|B)}^{\text{posterior}} = \frac{P(A, B)}{P(B)} = \overbrace{P(A)}^{\text{prior}} \times \frac{P(B|A)}{P(B)}$$

- ▶ *often counterintuitive*

(Monty Hall problem, gambling,  
human inability to generate random numbers, ...)

I have just thrown 10 successive sixes!

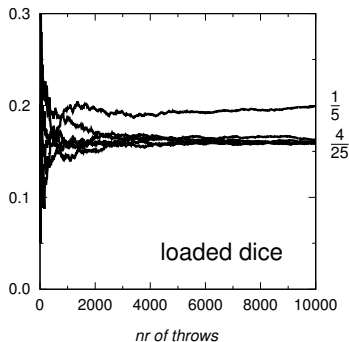
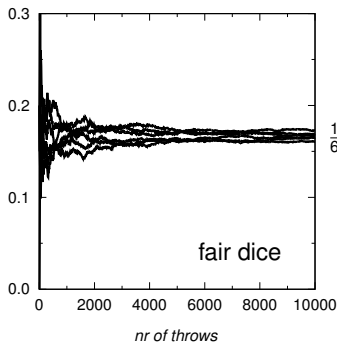
$Prob \approx 16.5 \cdot 10^{-8}$

*how likely am I to throw yet another six?*

- *how many data do we need to be sure of something?*

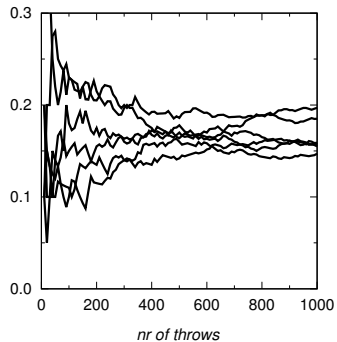
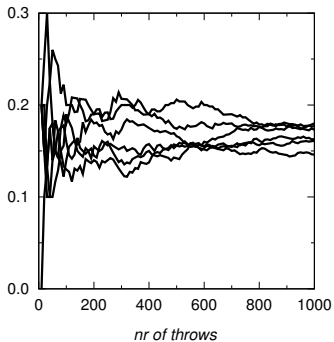
Is a genetic mutation harmless, or dangerous?

Is a given dice fair, or loaded?



typical data sets in cancer research:  
 $n \approx 500$  patients, at 2K£ each ...

what can we say  
after 500 throws?



- ▶ ‘probability’ can mean different things ...

our ignorance of

- (a) something *that cannot be known*  
(Russian roulette, we will spin the cylinder)
- (b) something *that is known, but not by us*  
(Russian roulette, cylinder has already been spun)

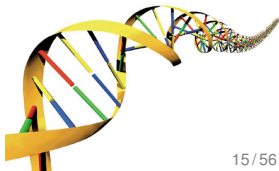


relevant in medicine?

suppose we find survival function  $S(t) = e^{-t/\tau}$

explanation I: homogeneous cohort, *random* death times,  
each individual  $i$  has hazard rate  $1/\tau$

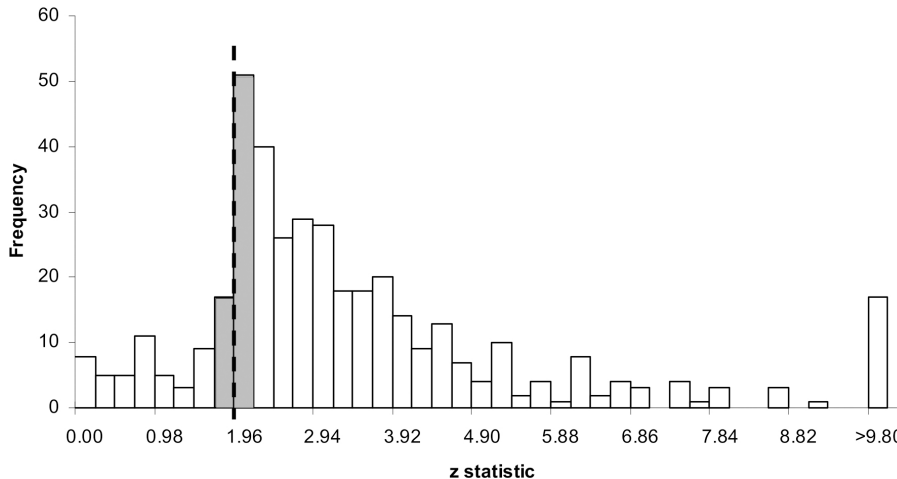
explanation II: heterogeneous cohort, *deterministic* death times  $t_i$ ,  
distributed according to  $p(t) = \tau^{-1}e^{-t/\tau}$   
(potential for stratification!)



## **z-scores**

reported in PLoS Medicine

Selective reporting  
(aka cheating)





Results from half of all clinical trials are hidden.  
Doctors don't have full information  
about the medicines we use.

Sign the petition



Donate >

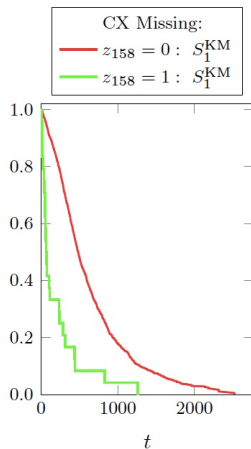
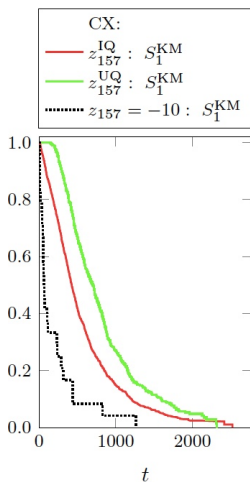


Get involved >



Latest news >

- ▶ *missing values in data sets. ...  
red herrings or white sharks?*



always check for  
informative missingness!

- correlation/association is not the same as causality!

imagine  $Z$   
is nr of hospital visits ...  
result: positive correlation  
between  $Z$  and risk!

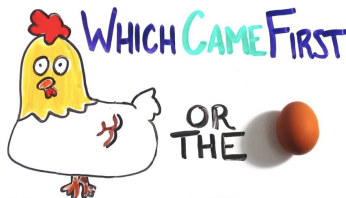
$\beta > 0$ , ergo hospital visits dangerous?

or:

$Z = 1, 0$ : given strong chemotherapy yes/no  
but treatment not given if patient too weak ...

result: positive association between  $Z$  and health!

protective effect reported, *even if chemotherapy ineffective!*



two chemotherapies, A and B,  
data on response rates from 880 patients

Q: which treatment is better?

	CHEMO A	CHEMO B
response rate	25% (76/300)	28% (162/580)

so treatment B is better,

now we zoom in ...

	CHEMO A	CHEMO B
medical centre 1	40% (40/100)	30% (150/500)
medical centre 2	18% (36/200)	15% (12/80)
response rate	25% (76/300)	28% (162/580)

still sure that B is better?

*Simpson's paradox*

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

### Terminology and objectives

- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Terminology and objectives

► *Data available*  $\mathcal{D} = \{(\mathbf{z}_1, r_1, t_1), \dots, (\mathbf{z}_N, r_N, t_N)\}$

$N$  samples/individuals  $(\mathbf{z}_i, r_i, t_i)$ ,  
drawn independently from  $p(t, r, \mathbf{z})$  (the population)

the ‘covariates’:

$\mathbf{z} \in \mathbb{R}^p$  :  *$p$  characteristics, measured at  $t=0$*

*uncontrolled* : e.g. gender, genome, ...

*controlled* : e.g. medical treatment, ...

*modifiable* : e.g. smoking, BMI, nutrition, ...

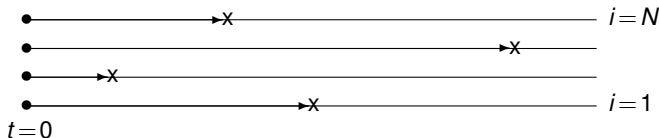
the ‘clinical outcome’:

$t \in \mathbb{R}^+$  : *failure time*  
*e.g. death, onset/recurrence of disease, ...*

$r \in \{0, 1, \dots, R\}$  : *risk type that triggered failure*  
 $r = 1 \dots R$  : *true risks/diseases*  
 $r = 0$  : *end of observation (‘censoring’)*

## Sheet1

pat	BMI	SELENIUM	PHYS_ACT_LEIS	PHYS_ACT_WORK	Smoking	Time	PCcens
1	22.63671875	105	2	1	0	33.6892539357	0
2	34.21875	65	2	2	0	24.810403833	1
3	20.06640625	72	2	3	2	23.1047227926	2
4	28.3984375	81	2	0	2	33.6783025325	0
5	22.94921875	73	0	3	2	32.8843258042	0
6	20.59765625	73	2	0	2	23.2936344969	1
7	26.46875	70	2	3	2	27.9069130732	2
8	26.38671875	91	1	0	1	26.6119096509	2
9	24.296875	73	-10	0	0	26.379192334	1
10	30.01953125	84	1	1	2	26.8254620123	2
11	25.4296875	95	1	0	2	30.6557152635	2
13	23.19921875	67	3	1	0	33.2457221081	0
14	22.90625	82	2	2	0	33.6399726215	0
15	21.62890625	53	1	1	1	33.2457221081	0
16	23.046875	77	2	0	2	33.6125941136	0
17	21.70703125	76	1	2	2	27.8466803559	1
18	22.91796875	102	1	0	2	33.6125941136	0
19	24.5078125	57	2	0	2	25.8097193703	2
20	26.58984375	72	1	2	2	26.803559206	2
21	26.76953125	78	2	0	2	33.234770705	0
22	20.4296875	75	-10	1	1	33.5934291581	0
23	25.0078125	69	0	3	0	33.6125941136	0
24	24.296875	73	2	3	2	21.6098562628	1
25	23.65625	75	2	3	1	33.2320328542	0
26	25.9296875	90	1	1	2	33.5359342916	0
27	23.3671875	58	1	1	2	33.2320328542	0
28	30.08984375	77	-10	-10	2	32.0438056126	2
29	31.08984375	66	1	0	0	29.4893908282	1
30	27.13671875	82	1	-10	1	33.1526351814	0
31	19.828125	68	2	2	2	28.2600958248	2
32	27.30859375	97	2	3	1	33.5742642026	0
33	23.41796875	77	2	0	0	25.9219712526	2



► *Objective*

find and quantify patterns that relate  
covariates  $\mathbf{z}$  to clinical outcomes  $(r, t)$ , in order to:

- *predict clinical outcome for individuals*
- *discover disease mechanisms*
- *design interventions (modifiable covariates)*

► *Complications*

- 'noise' caused by censoring
- we only know the earliest event  
(different risks prevent each other from happening)
- correlations between risks
- heterogeneity in patient cohorts
- overfitting danger, when  $p$  is large relative to  $N$



## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates**
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Survival probability and cause-specific hazard rates

### ► *Joint event time statistics*

imaginary situation: all events can be observed,

$t_r$ : time at which event  $r$  occurs,

event time distribution:  $P(t_0, \dots, t_R)$

normalisation: 
$$\int_0^\infty \dots \int_0^\infty dt_0 \dots dt_R P(t_0, \dots, t_R) = 1$$

### ► *Integrated event time distribution*

$$S(t_0, \dots, t_R) = \int_{t_0}^\infty \dots \int_{t_R}^\infty ds_0 \dots ds_R P(s_0, \dots, s_R)$$

meaning: probability that event 0 occurs later than  $t_0$ ,

and event 1 occurs later than  $t_1$ , and ...

$$S(0, \dots, 0) = \int_0^\infty \dots \int_0^\infty ds_0 \dots ds_R P(s_0, \dots, s_R) = 1$$

### ► *Survival function $S(t)$*

probability that *all* events happen later than time  $t$ :

$$S(t) = \int_t^\infty \dots \int_t^\infty ds_0 \dots ds_R P(s_0, \dots, s_R) = S(t, t, \dots, t)$$

► Cause-specific hazard rates  $h_r(t)$

how do individual risks impact on survival?

$$h_r(t) = - \left[ \frac{\partial}{\partial t_r} \log S(t_0, \dots, t_R) \right]_{t_k=t \text{ for all } k}$$

work out, using  $\frac{d}{dz}\theta(z) = \delta(z)$ :

$$\begin{aligned} h_r(t) &= \left[ \frac{\frac{\partial}{\partial t_r} \int_0^\infty \dots \int_0^\infty d\mathbf{s}_0 \dots d\mathbf{s}_R P(\mathbf{s}_0, \dots, \mathbf{s}_R) \prod_k \theta(s_k - t_k)}{S(t_0, \dots, t_R)} \right]_{t_k=t \forall k} \\ &= \left[ \frac{\int_0^\infty \dots \int_0^\infty d\mathbf{s}_0 \dots d\mathbf{s}_R P(\mathbf{s}_0, \dots, \mathbf{s}_R) \delta(s_r - t_r) \prod_{k \neq r} \theta(s_k - t_k)}{S(t_0, \dots, t_R)} \right]_{t_k=t \forall k} \\ &= \frac{\int_t^\infty \dots \int_t^\infty \left( \prod_{r \neq \mu}^R d\mathbf{s}_r \right) P(\mathbf{s}_0, \dots, \mathbf{s}_{\mu-1}, t, \mathbf{s}_{\mu+1}, \dots, \mathbf{s}_R)}{S(t)} \end{aligned}$$

$h_r(t)dt$ : probability that event  $r$  happens in time interval  $[t, t + dt)$ ,  
given that no event has happened yet prior to  $t$

$$h_r(t)dt = \text{Prob}\left(t_r \in [t, t+dt) \mid \text{no events yet at time } t\right) \quad (dt \downarrow 0)$$

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates**
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Data likelihood in terms of cause specific hazard rates

most of the relevant quantities in survival analysis  
can be written in terms of the cause specific hazard rates

### ► *Survival function*

$$\begin{aligned}\frac{d}{dt} \log S(t) &= \frac{d}{dt} \log S(t, t, \dots, t) \\ &= \sum_{r=0}^R \left[ \frac{\partial}{\partial t_r} \log S(t_0, \dots, t_R) \right]_{t_r=t \ \forall r} \\ &= - \sum_{r=0}^R h_r(t)\end{aligned}$$

Hence, using  $S(0) = 1$ ,

$$\log S(t) = \log S(0) - \int_0^t ds \sum_{r=0}^R h_r(s) = - \sum_{r=0}^R \int_0^t ds h_r(s)$$

result:

$$S(t) = e^{-\sum_{r=0}^R \int_0^t ds h_r(s)}$$

► *Data likelihood*

$p(t, r)dt$ : likelihood to observe *first* event being of type  $r$ ,  
and occurring in time interval  $[t, t + dt)$  (with  $dt \downarrow 0$ )

To observe the above, the following statements must be true:

time of the event is in  $[t, t + dt)$ ,

type of the event is  $r$ ,

no events occurred prior to  $t$

$$\theta(t_r - t)\theta(t + dt - t_r) \prod_{k \neq r} \theta(t_k - t) = 1$$

hence

$$\begin{aligned} p(t, r) &= \lim_{dt \downarrow 0} \frac{1}{dt} \text{Prob} \left( \theta(t_r - t)\theta(t + dt - t_r) \prod_{k \neq r} \theta(t_k - t) = 1 \right) \\ &= \lim_{dt \downarrow 0} \frac{1}{dt} \int_0^\infty \dots \int_0^\infty dt_0 \dots t_R P(t_1, \dots, t_R) \theta(t_r - t)\theta(t + dt - t_r) \prod_{k \neq r} \theta(t_k - t) \\ &= \int_0^\infty \dots \int_0^\infty dt_0 \dots t_R P(t_1, \dots, t_R) \lim_{\epsilon \downarrow 0} h_\epsilon(t_r - t) \prod_{k \neq r} \theta(t_k - t) \end{aligned}$$

$$h_\epsilon(z) = \epsilon^{-1} \theta(z) \theta(\epsilon - z) = \begin{cases} \epsilon^{-1} & \text{for } z \in [0, \epsilon] \\ 0 & \text{elsewhere} \end{cases}$$

note:  $\lim_{\epsilon \downarrow 0} h_\epsilon(z) = \delta(z)$ , so

$$\begin{aligned}
 p(t, r) &= \int_0^\infty \dots \int_0^\infty dt_0 \dots t_R P(t_1, \dots, t_R) \lim_{\epsilon \downarrow 0} h_\epsilon(t_r - t) \prod_{k \neq r} \theta(t_k - t) \\
 &= \int_0^\infty \dots \int_0^\infty dt_0 \dots t_R P(t_1, \dots, t_R) \delta(t_r - t) \prod_{k \neq r} \theta(t_k - t) \\
 &= h_t(t) S(t) = h_t(t) e^{-\sum_{r'=0}^R \int_0^t ds h_{r'}(s)}
 \end{aligned}$$

► *Further relation*

$$\begin{aligned}
 p(t) &= \sum_{r=0}^R p(t, r) = \left( \sum_{r=0}^R h_t(t) \right) e^{-\sum_{r'=0}^R \int_0^t ds h_{r'}(s)} \\
 &= -\frac{d}{dt} e^{-\sum_{r'=0}^R \int_0^t ds h_{r'}(s)} = -\frac{d}{dt} S(t)
 \end{aligned}$$

► *Cause-specific hazard rates in terms of data probabilities*

$$S(t) = S(0) + \int_0^t dt' \frac{d}{dt'} S(t') = 1 - \int_0^t dt' p(t') = \int_t^\infty ds p(s)$$

substitute into

formula for  $p(t, r)$ :

$$h_r(t) = \frac{p(t, r)}{\sum_{r'=0}^R \int_t^\infty ds p(s, r')}$$

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates

### **Pitfalls and misconceptions**

- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values



## Pitfalls and misconceptions

cause specific hazard rates can be tricky ...

$$S(t) = \prod_r e^{-\int_0^t ds \, h_r(s)}$$

- *Survival function formula factorizes over risks, does this imply that the risks are uncorrelated?*

No. All risks  $k \neq r$  can contribute to each  $h_r(t)$  via the conditioning, i.e. the likelihood that nothing has happened yet prior to  $t$ . Risks may well interact strongly with each other, but we can no longer see this after we have calculated the rates  $\{h_r(t)\}$  and forget about the times  $(t_0, \dots, t_R)$ .

- *Do we get the survival function for the situation where risk  $\mu$  is disabled (e.g. a disease removed from the world) by setting  $h_\mu(t)$  to zero?*

$$S(t) \rightarrow e^{-\sum_{r \neq \mu} \int_0^t ds \, h_r(s)}$$

No. We would have  $h_\mu(t) = 0$  for all  $t$ , but that is not all. Disabling risk  $\mu$  can change also *all* hazard rates  $h_r(t)$  with  $r \neq \mu$ , due to correlations among the different risks in combination with the conditioning.

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases**

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Special cases

- *Time-independent hazard rates*

$$h_r(t) = h_r: \quad S(t) = e^{-t \sum_{r=0}^R h_r} \quad p(t, r) = h_r e^{-t \sum_{r'=0}^R h_{r'}}$$

- *A single risk,  $R=1$*

One risk,  
hazard rate  $h(t)$ :  $S(t) = e^{-\int_0^t ds h(s)} \quad p(t) = h(t) e^{-\int_0^t ds h(s)}$

- *Most probable event time distribution for  $R=1$*

Suppose we know only average event time  $\langle t \rangle = \int_0^\infty dt t p(t)$ ,  
most probable  $p(t)$ : maximize Shannon entropy  
 $H = -\int_0^\infty dt p(t) \log p(t)$ , subject to  $\int_0^\infty dt p(t) = 1$  and  $\int_0^\infty dt t p(t) = \langle t \rangle$

Lagrange's method:

$$\frac{\delta}{\delta p(t)} \int_0^\infty ds p(s) \log p(s) = \frac{\delta}{\delta p(t)} \left\{ \lambda_0 \int_0^\infty ds p(s) + \lambda_1 \int_0^\infty ds p(s) s \right\}$$

$$1 + \log p(t) = \lambda_0 + \lambda_1 t \quad \text{so} \quad p(t) = e^{\lambda_0 - 1 + \lambda_1 t}$$

use constraints:

$$p(t) = \langle t \rangle^{-1} e^{-t/\langle t \rangle}$$

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times**
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Independently distributed event times

$$P(t_0, \dots, t_R) = \prod_{r=0}^R P_r(t_r)$$

### ► Integrated event time distr

$$S(t_0, \dots, t_R) = \int_{t_0}^{\infty} \dots \int_{t_R}^{\infty} ds_0 \dots ds_R \prod_{r=0}^R P_r(t_r) = \prod_{r=0}^R S_r(t_r)$$
$$S_r(t) = \int_t^{\infty} ds P_r(s)$$

### ► Cause specific hazard rates

$$h_r(t) = - \left[ \frac{\partial}{\partial t_r} \sum_{r'=0}^R \log S_{r'}(t_{r'}) \right]_{t_k=t \text{ for all } k} = - \frac{d}{dt} \log S_r(t)$$

hence

$$S_r(t) = e^{-\int_0^t ds h_r(s)}$$

joint event time distr now follows from the cause-specific hazard rates

$$P_r(t) = - \frac{d}{dt} S_r(t) = - \frac{d}{dt} e^{-\int_0^t ds h_r(s)} = h_t(t) e^{-\int_0^t ds h_r(s)}$$

$$P(t_0, \dots, t_R) = \prod_{r=0}^R \left[ h_t(t_r) e^{-\int_0^{t_r} ds h_r(s)} \right]$$

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## The identifiability problem (Tsiatis)

- *Observable from data:*  $p(t, r)$ ,  
equivalently:  $\{h_0(t), \dots, h_R(t)\}$ , since

$$p(t, r) = h_r(t) e^{-\sum_{r'=0}^R \int_0^t ds h_{r'}(s)}, \quad h_r(t) = \frac{p(t, r)}{\sum_{r'=0}^R \int_t^\infty ds p(s, r')}$$

- For any  $\{h_0(t), \dots, h_R(t)\}$ , even those corresponding to *statistically dependent event times*, there exists a distribution for *independent* event times that will give exactly the same cause-specific hazard rates, namely

$$P(t_0, \dots, t_R) = \prod_{r=0}^R \prod_{r=0}^R \left[ h_r(t_r) e^{-\int_0^{t_r} ds h_r(s)} \right]$$

Hence, survival data alone do not generally permit us to identify the underlying joint distribution of event times – in particular, we cannot infer whether or not the event times of the different risks are independent.

a serious problem ...

- ▶ *The Bayesian view*  
on the identifiability problem
  - multiple hypotheses  $H$  may explain our data
  - but not all are equally probable ...
  - calculate each  $\text{Prob}(H|\mathcal{D})$  from Bayes' formula

▶ *Illustration*

true data:

$$p(t_2) = ae^{-at_2}, \quad \begin{cases} \text{with prob } \epsilon: & t_1 = t_2 + \tau \\ \text{with prob } 1-\epsilon: & \text{draw } t_1 \text{ from } p(t_1) = be^{-bt_1} \end{cases}$$

explanation assuming  
risk independence:

$$p(t_2) = ae^{-at_2}, \quad p(t_1) = - \underbrace{\left( \epsilon + (1-\epsilon)e^{-bt_1} \right) \log \left( \epsilon + (1-\epsilon)e^{-bt_1} \right)}_{\text{with prob } \epsilon: \text{ event 1 never happens}}$$

implausible if e.g. risk 2 is cancer, risk 1 is death ...



## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions**
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Proportional hazards regression (Cox) definitions and assumptions

### ► *Survival analysis with covariates*

Add  $\mathbf{z}$  as conditions to definitions and identities

$$S(t) \rightarrow S(t|\mathbf{z}), \quad h_r(t) \rightarrow h_r(t|\mathbf{z}), \quad p(t, r) \rightarrow p(t, r|\mathbf{z})$$

$$S(t|\mathbf{z}) = \int_t^\infty \dots \int_t^\infty ds_0 \dots ds_R P(s_0, \dots, s_R|\mathbf{z})$$

$$S(t|\mathbf{z}) = e^{-\sum_{r=0}^R \int_0^t ds \, h_r(s|\mathbf{z})}, \quad p(t, r|\mathbf{z}) = h_r(t|\mathbf{z}) e^{-\sum_{r'=0}^R \int_0^t ds \, h_{r'}(s|\mathbf{z})}$$

### ► *Cox model*

Parametrized form for the hazard rates:

$$h_r(t|\mathbf{Z}) = \lambda_r(t) e^{\beta^r \cdot \mathbf{z}}, \quad \beta^r = (\beta_1^r, \dots, \beta_p^r), \quad \beta^r \cdot \mathbf{z} = \sum_{\mu=1}^p \beta_\mu^r z_\mu$$

$\lambda_r(t)$ : ‘base hazard rate’ of risk  $r$   
(covariate-independent contribution to risk)

$\beta^r$ : ‘association parameters’ of risk  $r$   
(impact of covariate values on risk)

► *Main features of Cox's choice*

► *'Proportional hazards'*

due to exponential form, effect of each covariate is multiplicative:

$$h_r(t) = \underbrace{\lambda_r(t)}_{\text{base hazard rate}} \times \underbrace{e^{\beta_1^r z_1} \times \dots \times e^{\beta_p^r z_p}}_{\text{'proportional hazards'}}$$

► *Effects of covariates on risk independent of time*

► *There exists a hyper-plane in covariate space that separates high risk individuals from low risk individuals*

*'high risk individuals'* :  $\beta_1^r z_1 + \dots + \beta_p^r z_p$  *large*

*'low risk individuals'*  $\beta_1^r z_1 + \dots + \beta_p^r z_p$  *small*

► *One can quantify risk impact of each individual covariate  $\mu$  in a single time-independent number: the 'hazard ratio'*

$$HR_{\mu}^r = \frac{h_r(t|\mathbf{z})|_{z_{\mu}=1}}{h_r(t|\mathbf{z})|_{z_{\mu}=0}} = \frac{\lambda_r(t)e^{\beta_{\mu}^r \cdot 1 + \sum_{\nu \neq \mu} \beta_{\nu}^r z_{\nu}}}{\lambda_r(t)e^{\beta_{\mu}^r \cdot 0 + \sum_{\nu \neq \mu} \beta_{\nu}^r z_{\nu}}} = e^{\beta_{\mu}^r}$$

If no impact on risk:  $\beta_{\mu}^r = 0$ ,  $HR_{\mu}^r = 1$ .

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data**
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Detour: parameter estimation from data

given data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  
and a model  $p(\mathbf{x}|\theta)$  to explain these data,  
what can we say about the parameters  $\theta$ ?

### ► Bayesian parameter inference

assume the  $\{\mathbf{x}_i\}$  were indeed drawn  
randomly & independently from  $p(\mathbf{x}|\theta)$ ,

$$p(\mathcal{D}|\theta) = \prod_{i=1}^N p(\mathbf{x}_i|\theta)$$

Bayes' identity:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int d\theta' p(\mathcal{D}|\theta')p(\theta')}, \quad p(\theta): \text{prior}$$

### ► Simplifications

$$MAP: \quad \theta^* = \operatorname{argmax} p(\theta|\mathcal{D}) = \operatorname{argmin} \left[ \overbrace{-\log p(\mathcal{D}|\theta)}^{\text{minus log-likelihood}} \overbrace{-\log p(\theta)}^{\text{regularizer}} \right]$$

$$ML: \quad \theta^* = \operatorname{argmin} \left[ -\log p(\mathcal{D}|\theta) \right] \quad \text{i.e. } p(\theta) = \text{constant}$$

► *Maximum Likelihood (ML) regression*

$$\boldsymbol{\theta}^* = \operatorname{argmin}[-\log p(\mathcal{D}|\boldsymbol{\theta})]$$

define empirical  
data distribution

$$\hat{p}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}_i)$$

note:

$$\begin{aligned} -\frac{1}{N} \log p(\mathcal{D}|\boldsymbol{\theta}) &= -\frac{1}{N} \log \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N \log p(\mathbf{x}_i|\boldsymbol{\theta}) \\ &= -\int d\mathbf{x} \hat{p}(\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) \\ &= \int d\mathbf{x} \hat{p}(\mathbf{x}) \log \left[ \frac{\hat{p}(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} \right] - \int d\mathbf{x} \hat{p}(\mathbf{x}) \log \hat{p}(\mathbf{x}) \\ &= \underbrace{D(\hat{p}||p_{\boldsymbol{\theta}})}_{\text{KL distance}} + \underbrace{H[\hat{p}]}_{\text{Shannon entropy}} \end{aligned}$$

hence: ML finds the parameter vector  $\boldsymbol{\theta}$  that minimizes  
the KL distance between  $\hat{p}(\mathbf{x})$  and  $p(\mathbf{x}|\boldsymbol{\theta})$

- *Beyond most probable parameters:  
error bars*

return to full posterior distribution

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{e^{-\Omega(\boldsymbol{\theta}, \mathcal{D})}}{\int d\boldsymbol{\theta}' e^{-\Omega(\boldsymbol{\theta}', \mathcal{D})}}, \quad \Omega(\boldsymbol{\theta}, \mathcal{D}) = -\log p(\mathcal{D}|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta})$$

expand  $\Omega$  around minimum  $\boldsymbol{\theta}^*$ :

$$\Omega(\boldsymbol{\theta}, \mathcal{D}) = \Omega(\boldsymbol{\theta}^*, \mathcal{D}) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \cdot \mathbf{A}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \mathcal{O}(|\boldsymbol{\theta} - \boldsymbol{\theta}^*|^3)$$

truncate after quadratic term:

$$p(\boldsymbol{\theta}|\mathcal{D}) \approx \left[ \frac{\det \mathbf{A}}{(2\pi)^N} \right]^{\frac{1}{2}} e^{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \cdot \mathbf{A}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)}, \quad \langle (\theta_\mu - \theta_\mu^*)(\theta_\nu - \theta_\nu^*) \rangle = (\mathbf{A}^{-1})_{\mu\nu}$$

hence, error bars for  
MAP/ML estimators:

$$\theta_\mu = \theta_\mu^* \pm (\mathbf{A}^{-1})_{\mu\mu}, \quad A_{\mu\nu} = \frac{\partial^2}{\partial \theta_\mu \partial \theta_\nu} \left[ -\log p(\mathcal{D}|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) \right]_{\boldsymbol{\theta}^*}$$

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model**
- Uniqueness, error bars, and  $p$ -values



## ML parameters of Cox's model

$$\boldsymbol{\theta} = \{\lambda_r, \boldsymbol{\beta}^r\} :$$

$$p(t, r | \mathbf{z}, \boldsymbol{\theta}) = h_r(t | \mathbf{z}, \boldsymbol{\theta}) e^{-\sum_{r'=0}^R \int_0^t ds h_{r'}(s | \mathbf{z}, \boldsymbol{\theta})}, \quad h_r(t | \mathbf{z}, \boldsymbol{\theta}) = \lambda_r(t) e^{\boldsymbol{\beta}^r \cdot \mathbf{z}}$$

### ► ML inference

$$\begin{aligned} \boldsymbol{\theta}^* &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ -\log p(\mathcal{D} | \boldsymbol{\theta}) \right] \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ -\sum_{i=1}^N \log p(t_i, r_i | \mathbf{z}_i, \boldsymbol{\theta}) \right] \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ -\sum_{i=1}^N \log h_{r_i}(t_i | \mathbf{z}_i, \boldsymbol{\theta}) + \sum_{i=1}^N \sum_{r=0}^R \int_0^{t_i} dt \lambda_r(t) e^{\boldsymbol{\beta}^r \cdot \mathbf{z}_i} \right] \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ -\sum_{i=1}^N \log \lambda_{r_i}(t_i) - \sum_{i=1}^N \boldsymbol{\beta}^{r_i} \cdot \mathbf{z}_i + \sum_{i=1}^N \sum_{r=0}^R \int_0^{t_i} dt \lambda_r(t) e^{\boldsymbol{\beta}^r \cdot \mathbf{z}_i} \right] \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{r=0}^R \left[ -\sum_{i=1}^N \delta_{r, r_i} \log \lambda_r(t_i) - \sum_{i=1}^N \delta_{r, r_i} \boldsymbol{\beta}^r \cdot \mathbf{z}_i + \sum_{i=1}^N \int_0^{t_i} dt \lambda_r(t) e^{\boldsymbol{\beta}^r \cdot \mathbf{z}_i} \right] \end{aligned}$$

To minimize

$$\begin{aligned}\Psi[\{\lambda_r, \beta^r\}] &= \sum_{r=0}^R \left[ - \int dt \log \lambda_r(t) \sum_{i=1}^N \delta_{r,r_i} \delta(t-t_i) - \beta^r \cdot \sum_{i=1}^N \delta_{r,r_i} \mathbf{z}_i \right. \\ &\quad \left. + \int dt \lambda_r(t) \sum_{i=1}^N \theta(t_i-t) e^{\beta^r \cdot \mathbf{z}_i} \right]\end{aligned}$$

► Minimize over functions  $\lambda_r(t)$  first

$$\frac{\delta}{\delta \lambda_r(t)} \left[ \sum_{i=1}^N \int dt \theta(t_i-t) \lambda_r(t) e^{\beta^r \cdot \mathbf{z}_i} - \int dt \log \lambda_r(t) \sum_{i=1}^N \delta_{r,r_i} \delta(t-t_i) \right] = 0$$

(functional differentiation)

$$\sum_{i=1}^N \theta(t_i-t) e^{\beta^r \cdot \mathbf{z}_i} - \frac{1}{\lambda_r(t)} \sum_{i=1}^N \delta_{r,r_i} \delta(t-t_i) = 0$$

$$\lambda_r(t) = \frac{\sum_{i=1}^N \delta_{r,r_i} \delta(t-t_i)}{\sum_{i=1}^N \theta(t_i-t) e^{\beta^r \cdot \mathbf{z}_i}} \quad \text{Breslow's estimator}$$

► Insert into  $\Psi$ , then minimize over  $\{\beta^r\}$

► To minimize

$$\begin{aligned}
 \Psi[\{\beta^r\}] &= \sum_{r=0}^R \left[ - \int dt \log \left( \frac{\sum_{i=1}^N \delta_{r,r_i} \delta(t-t_i)}{\sum_{i=1}^N \theta(t_i-t) e^{\beta^r \cdot \mathbf{z}_i}} \right) \sum_{j=1}^N \delta_{r,r_j} \delta(t-t_j) \right. \\
 &\quad \left. - \beta^r \cdot \sum_{i=1}^N \delta_{r,r_i} \mathbf{z}_i + \int dt \sum_{i=1}^N \delta_{r,r_i} \delta(t-t_i) \right] \\
 &= \sum_{r=0}^R \left[ \int dt \log \left( \sum_{i=1}^N \theta(t_i-t) e^{\beta^r \cdot \mathbf{z}_i} \right) \sum_{j=1}^N \delta_{r,r_j} \delta(t-t_j) \right. \\
 &\quad \left. - \beta^r \cdot \sum_{i=1}^N \delta_{r,r_i} \mathbf{z}_i \right] + \text{terms independent of } \{\beta^r\} \\
 &= \sum_{r=0}^R \Psi_r(\beta^r) + \text{terms independent of } \{\beta^r\}
 \end{aligned}$$

with

$$\Psi_r(\beta) = \sum_{j=1}^N \delta_{r,r_j} \log \left( \sum_{i=1}^N \theta(t_i-t_j) e^{\beta \cdot \mathbf{z}_i} \right) - \beta \cdot \sum_{i=1}^N \delta_{r,r_i} \mathbf{z}_i$$

hence

$$\beta^{r*} = \operatorname{argmin}_{\beta} \Psi_r(\beta)$$

- Each risk  $r$ , find minima of  $\Psi_r$ :

$$\begin{aligned}\frac{\partial}{\partial \beta_\mu} \Psi_r(\beta) &= \frac{\partial}{\partial \beta_\mu} \left[ \sum_{j=1}^N \delta_{r,r_j} \log \left( \sum_{i=1}^N \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i} \right) - \beta \cdot \sum_{j=1}^N \delta_{r,r_j} \mathbf{z}_j \right] \\ &= \sum_{j=1}^N \delta_{r,r_j} \frac{\sum_{i=1}^N z_{i\mu} \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}{\sum_{i=1}^N \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}} - \sum_{j=1}^N \delta_{r,r_j} z_{j\mu}\end{aligned}$$

so  $\beta^{r*}$  is solution of:

$$\sum_{j=1}^N \delta_{r,r_j} \left[ \frac{\sum_{i=1}^N z_{i\mu} \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}{\sum_{i=1}^N \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}} - z_{j\mu} \right] = 0$$

$p$  coupled nonlinear equations, for each risk

Final protocol:

1. Solve  $\{\beta^{r*}\}$  from above eqns (numerically)
2. Calculate  $\{\lambda^r(t)\}$  (from Breslow's formula)
3. Predict outcomes via  $p(r, t|\mathbf{z})$  for new patients (using Cox's model, with ML parameters)

## Introduction

- Post-genome medicine
- Statistics in medicine is tricky

## The formalism of survival analysis

- Terminology and objectives
- Survival probability and cause-specific hazard rates
- Data likelihood in terms of cause specific hazard rates
- Pitfalls and misconceptions
- Special cases

## Event time correlations and identifiability

- Independently distributed event times
- The identifiability problem

## Proportional hazards regression

- Definitions and assumptions
- Parameter estimation from data
- ML parameters of Cox's model
- Uniqueness, error bars, and  $p$ -values

## Uniqueness, error bars, and and $p$ -values

Curvature of  $\Psi_r(\beta)$ :

$$\begin{aligned}\frac{\partial^2 \Psi_r(\beta)}{\partial \beta_\mu \partial \beta_\nu} &= \frac{\partial}{\partial \beta_\nu} \left[ \sum_{j=1}^N \delta_{r,r_j} \frac{\sum_{i=1}^N z_{i\mu} \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}{\sum_{i=1}^N \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}} - \sum_{j=1}^N \delta_{r,r_j} z_{j\mu} \right] \\ &= \sum_{j=1}^N \delta_{r,r_j} \left[ \langle z_\mu z_\nu \rangle_j - \langle z_\mu \rangle_j \langle z_\nu \rangle_j \right]\end{aligned}$$

with

$$\langle f(\mathbf{z}) \rangle_j = \sum_{i=1}^N w(i|j) f(\mathbf{z}_i), \quad w(i|j) = \frac{\theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}{\sum_{i=1}^N \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}$$

properties, consequences:

- *Convexity*

curvature matrix is positive definite, i.e.  $\Psi_r(\beta)$  convex, since

$$\begin{aligned}\forall \mathbf{x} \in \mathbb{R}^p : \quad \sum_{\mu, \nu=1}^p x_\mu x_\nu \frac{\partial^2 \Psi_r(\beta)}{\partial \beta_\mu \partial \beta_\nu} &= \sum_{j=1}^N \delta_{r,r_j} \left[ \langle (\mathbf{x} \cdot \mathbf{z})^2 \rangle_j - \langle \mathbf{x} \cdot \mathbf{z} \rangle_j^2 \right] \\ &= \sum_{j=1}^N \delta_{r,r_j} \left[ \langle (\mathbf{x} \cdot \mathbf{z} - \langle \mathbf{x} \cdot \mathbf{z} \rangle_j)^2 \rangle_j \right] \geq 0\end{aligned}$$

- *Uniqueness of  $\beta^{f*}$*

since  $\Psi_r(\beta)$  convex, can have only one minimum

► *Error bars for association parameters*

Neglect fluctuations in  $\{\lambda_r(t)\}$ , focus on  $\Delta\beta_\mu^{r\star}$ :

$$A(r)_{\mu\nu} = \sum_{j=1}^N \delta_{r,r_j} \left[ \langle z_\mu z_\nu \rangle_j - \langle z_\mu \rangle_j \langle z_\nu \rangle_j \right]$$

$$\langle f(\mathbf{z}) \rangle_j = \sum_{i=1}^N w(i|j) f(\mathbf{z}_i), \quad w(i|j) = \frac{\theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}{\sum_{i=1}^N \theta(t_i - t_j) e^{\beta \cdot \mathbf{z}_i}}$$

Then

$$\beta_\mu^r = \beta_\mu^{r\star} \pm \sigma_\mu^r, \quad \sigma_\mu^r = (\mathbf{A}(r)^{-1})_{\mu\mu}$$

► *p-values of inferred  $\beta_\mu^\star$*

definition: the probability to find a value  $\beta_\mu^\star$  (or one further away from zero) due to fluctuations, when the true value is zero

approx: assume above error bar is correct, disregard correlations,

$$\begin{aligned} p\text{-value} &= \text{Prob}\left(|\beta_\mu| \geq |\beta_\mu^\star|\right) = 1 - \frac{2}{\sigma_\mu \sqrt{2\pi}} \int_0^{|\beta_\mu^\star|} d\beta \, e^{-\frac{1}{2}\beta^2/\sigma_\mu^2} \\ &= 1 - \text{Erf}\left(|\beta_\mu^\star|/\sigma_\mu \sqrt{2}\right), \quad \beta_\mu^\star/\sigma_\mu : \text{z-score} \end{aligned}$$

## Explanation for Simpson's paradox

	CHEMO A	CHEMO B
medical centre 1	40% (40/100)	30% (150/500)
medical centre 2	18% (36/200)	15% (12/80)
response rate	25% (76/300)	28% (162/580)

$$P(\text{response}|\text{chemo}) = \sum_{\text{centres}} P(\text{response}|\text{chemo}, \text{centre}) P(\text{centre}|\text{chemo})$$

$$P(\text{resp}|A) = \frac{40}{100} \frac{100}{300} + \frac{36}{200} \frac{200}{300} = 25\%$$

$$P(\text{resp}|B) = \frac{30}{100} \frac{500}{580} + \frac{15}{100} \frac{80}{580} = 28\%$$

if chemo choice indep of centre:

$$P(\text{resp}|A) = \frac{40}{100} \frac{1}{2} + \frac{36}{200} \frac{1}{2} = 29\%$$

$$P(\text{resp}|B) = \frac{30}{100} \frac{1}{2} + \frac{15}{100} \frac{1}{2} = 22.5\%$$