

Towards a theory of overfitting in proportional hazards regression for survival data

ACC Coolen

King's College London

April 19th 2016

Introduction

- Regression for time-to-event data

- Overfitting in PH regression

Replica analysis of overfitting in PH regression

- The basic ideas

- Translation to Cox's model

- Replica symmetric solution

- The end game

Simple approximations

Discussion

Introduction

Regression for time-to-event data

Overfitting in PH regression

Replica analysis of overfitting in PH regression

The basic ideas

Translation to Cox's model

Replica symmetric solution

The end game

Simple approximations

Discussion

Regression for time-to-event data

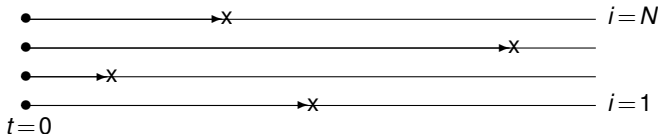
- *Data* $\mathcal{D} = \{(\mathbf{z}_1, t_1), \dots, (\mathbf{z}_N, t_N)\}$

samples (\mathbf{z}_i, t_i) ,

drawn indep from $p(t, \mathbf{z})$

$\mathbf{z}_i \in \mathbb{R}^d$: d covariates (measured at $t = 0$)

$t_i \in \mathbb{R}^+$: failure time (death, onset of disease, ...)



- *Objective*

find and quantify patterns that relate covariates to event times, in order to:

1. *predict clinical outcome for individuals*
2. *discover disease mechanisms*
3. *design interventions (modifiable covariates)*

Proportional hazards regression (DR Cox, 1972)



hazard rate : $h(t|\mathbf{z}) = \lambda(t)e^{\beta \cdot \mathbf{z}}$

event time dist : $p(t|\mathbf{z}, \beta, \lambda) = -\frac{d}{dt} \exp\left[-e^{\beta \cdot \mathbf{z}} \int_0^t dt' \lambda(t')\right]$

parameters : $\beta = (\beta_1, \dots, \beta_d), \quad \lambda(t) \quad t \geq 0$

- ▶ Maximum Likelihood estimation

$$(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \left\{ \frac{1}{N} \sum_i \log p(t_i | \mathbf{z}_i, \beta, \lambda) \right\}$$

- ▶ Maximise over $\lambda(t)$ *first*

$$\hat{\lambda}(t|\beta) = \frac{\sum_j \delta(t-t_j)}{\sum_k \theta(t_k-t) e^{\beta \cdot \mathbf{z}_k}} \quad (\text{Breslow estimator})$$

$$\hat{\beta} = \operatorname{argmax}_{\beta} \left\{ \sum_i \beta \cdot \mathbf{z}_i - \sum_i \log \left[\frac{\sum_j e^{\beta \cdot \mathbf{z}_j} \theta(t_j-t_i)}{\sum_j \theta(t_j-t_i)} \right] \right\}$$

relatively simple and computationally painless,
extremely successful,
still the main tool of medical statisticians ...



Beyond the basic model ...

► *Fine tuning*

- include left- right- or interval censoring (slightly different formula $p(\mathcal{D}|\beta, \lambda)$)
- consistent base hazard rate, such that $\int_0^\infty dt \lambda(t) = \infty$ (ML subject to constraint $\int_0^\infty dt \lambda(t) = R$, then $R \rightarrow \infty$)

► *Multiple risks*

risk labels $r_i \in \{0, \dots, R\}$,

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

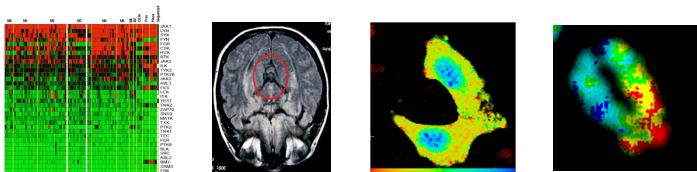
► *Frailty, random effects and latent class models*

(simple formulae only for special choices)

$$p(t|\mathbf{z}, \beta, \lambda) = -\frac{d}{dt} \sum_{\ell=1}^L w_\ell \exp\left[-e^{\beta^\ell \cdot \mathbf{z}} \int_0^t dt' \lambda^\ell(t')\right]$$

What has changed since the 1970s?

- ▶ *Medical data have evolved*



- ▶ *sheer volume ...*
- ▶ *diversity* of data sources
(clinical, genomic, biomarkers, health records, imaging, ...)
- ▶ *complexity* of experimental pipelines
(confounders, batch effects, variability between centres, ...)
- ▶ *dimension mismatch*
then: ~ 500 samples, ~ 10 covariates
now: ~ 1000 samples, $\sim 10^6$ covariates

Introduction

Regression for time-to-event data

Overfitting in PH regression

Replica analysis of overfitting in PH regression

The basic ideas

Translation to Cox's model

Replica symmetric solution

The end game

Simple approximations

Discussion

overfitting in Cox regression

ML method ...

p-values, z-scores,
confidence intervals
don't measure overfitting!

rule of thumb: $d_{\max} = N/10$

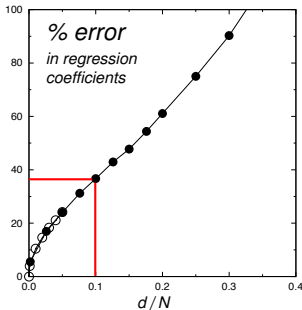
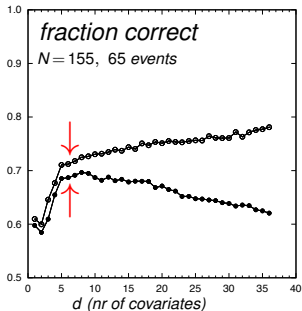
- ▶ too optimistic ...
- ▶ must depend on β ...
- ▶ must depend on cov correlations ...

*Analytical theory
of overfitting in
Cox regression?*

uncorrelated covars

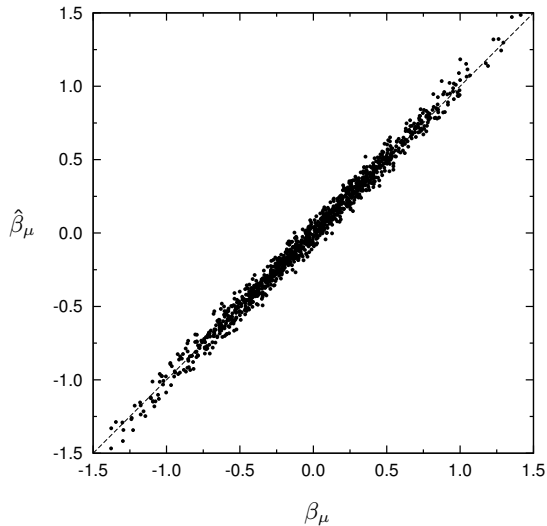
○: $N = 1000$

●: $N = 500$



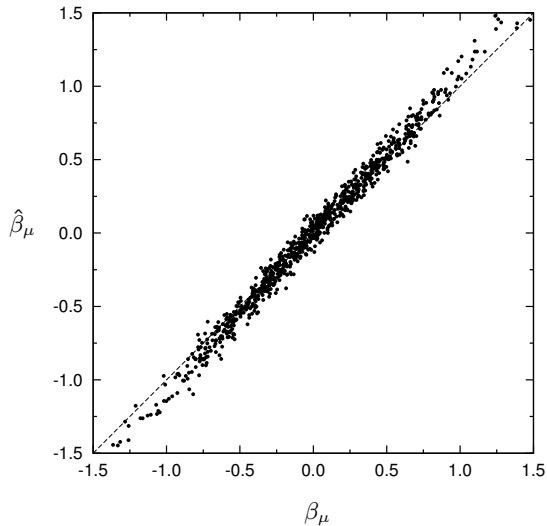
$N = 500$,
predicted versus true regression coefficients

$d/N = 0.002$



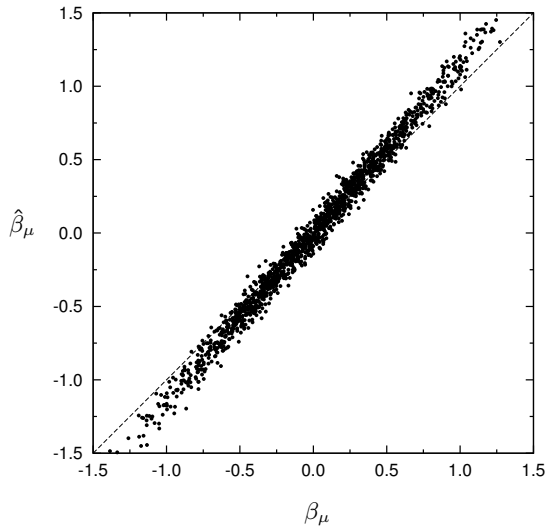
$N = 500$,
predicted versus true regression coefficients

$d/N = 0.10$



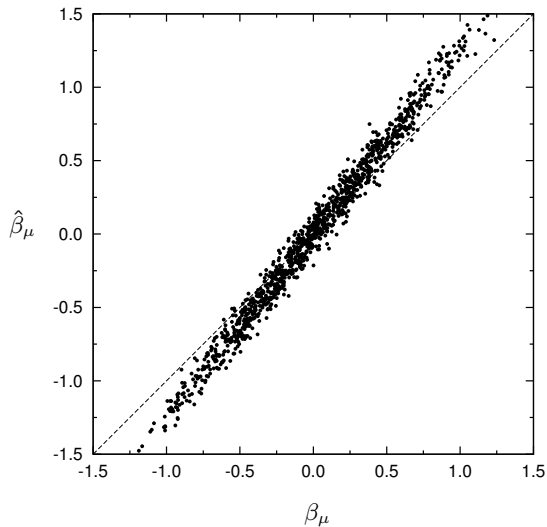
$N = 500$,
predicted versus true regression coefficients

$d/N = 0.20$



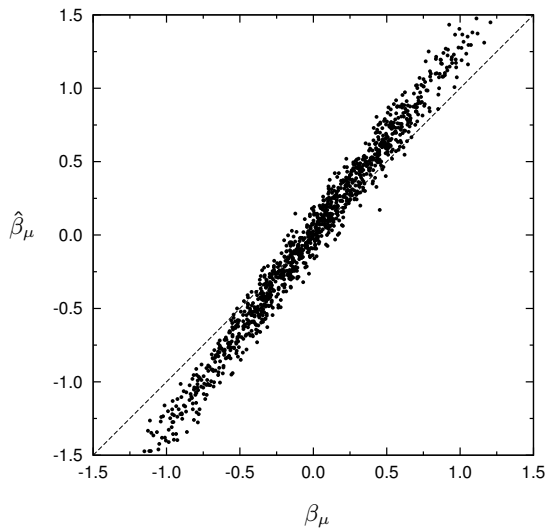
$N = 500$,
predicted versus true regression coefficients

$d/N = 0.30$



$N = 500$,
predicted versus true regression coefficients

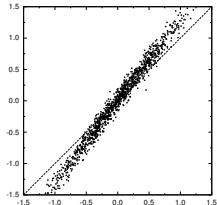
$d/N = 0.40$



Bad news

Overfitting *more dangerous*
than finite sample noise ...

*we always overestimate the
strength of associations
(whether positive or negative)*



Good news

Unlike pure noise,
deterministic bias may be predictable ...

New possibilities, roadmap for research ...

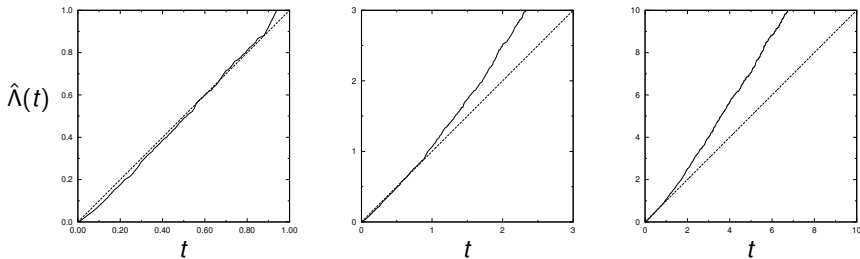
- ▶ Predict asymptotic impact of overfitting, in terms of
 - ratio d/N
 - correlations among covariates
 - true association strengths β
- ▶ Overfitting correction of Cox parameters
 - reliable regression at ratios $d/N \sim 0.5$?
 - Bayesian Cox regression with unbiased estimates?

$N=500$, $d/N=0.2$

inferred integrated base hazard rate

$$\hat{\Lambda}(t) = \int_0^t dt' \hat{\lambda}(t')$$

dashed: true values



- ▶ deviations increase with ratio $\alpha = d/N$
- ▶ here: synthetic data with $\lambda(t) = 1$,
so average event time for $\mathbf{z} = \mathbf{0}$ is 1

Intuition for the problem ...

- ▶ *Overfitting in ML regression*

assumed model: p_{θ}

$$\theta_{\text{ML}} = \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta) = \operatorname{argmin}_{\theta} D(\hat{p}||p_{\theta})$$

$$\hat{p}(t, \mathbf{z}) = \frac{1}{N} \sum_i \delta(t-t_i) \delta(\mathbf{z}-\mathbf{z}_i), \quad D(\hat{p}||p_{\theta}) = \int dt d\mathbf{z} \hat{p}(t, \mathbf{z}) \log \left[\frac{\hat{p}(t|\mathbf{z})}{p(t|\mathbf{z}, \theta)} \right]$$

ML regression: move $p(t|\mathbf{z}, \theta)$ towards $\hat{p}(t|\mathbf{z})$

true pars: θ^*

- ▶ fixed d : $\lim_{N \rightarrow \infty} \hat{p}(t, \mathbf{z}) = p(t, \mathbf{z}|\theta^*)$, so $\theta_{\text{ML}} = \theta^*$
 - ▶ $d = \mathcal{O}(N)$: $\lim_{N \rightarrow \infty} \hat{p}(t, \mathbf{z}) \neq p(t, \mathbf{z}|\theta^*)$...
-
- ▶ *Barrier to overfitting theory*

want: study relation between $\theta_{\text{ML}}(\mathcal{D})$ and θ^* , for $d = \mathcal{O}(N)$

need: formula for $\theta_{\text{ML}}(\mathcal{D})$...

Cox regression: $\theta_{\text{ML}}(\mathcal{D})$ solved from transcendental eqn

Introduction

Regression for time-to-event data

Overfitting in PH regression

Replica analysis of overfitting in PH regression

The basic ideas

Translation to Cox's model

Replica symmetric solution

The end game

Simple approximations

Discussion

The basic ideas

Step 1 – define a suitable overfitting measure

- ▶ Let \hat{p}_{θ^*} be empirical distr of (t, \mathbf{z}) ,
for data with true pars θ^*

note that

$$\theta_{\text{ML}} = \operatorname{argmin}_{\theta} D(\hat{p}_{\theta^*} \| p_{\theta})$$

$$\theta = \theta^*: D(\hat{p}_{\theta^*} \| p_{\theta}) = D(\hat{p}_{\theta^*} \| p_{\theta^*}) \leftarrow \text{not zero!}$$

Define:

$$E(\theta^*, \mathcal{D}) = \min_{\theta} D(\hat{p}_{\theta^*} \| p_{\theta}) - D(\hat{p}_{\theta^*} \| p_{\theta^*})$$

$E(\theta^*, \mathcal{D}) > 0$: underfitting

$E(\theta^*, \mathcal{D}) = 0$: optimal fitting

$E(\theta^*, \mathcal{D}) < 0$: overfitting

- ▶ *Typical behaviour*

$$\begin{aligned} E(\theta^*) &= \left\langle E(\theta^*, \mathcal{D}) \right\rangle_{\mathcal{D}} \\ &= \left\langle \min_{\theta} \left\{ \frac{1}{N} \sum_i \log \left[\frac{p(t_i | \mathbf{z}_i, \theta^*)}{p(t_i | \mathbf{z}_i, \theta)} \right] \right\} \right\rangle_{\mathcal{D}} \end{aligned}$$

□

Step 2 – eliminate minimisation over β

- ▶ *Laplace identity*
(steepest descent)

$$\lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \log \int dx e^{\gamma f(x)} = \max_x f(x)$$

use in reverse:

$$\begin{aligned} E(\theta^*) &= \left\langle \min_{\theta} \left\{ \frac{1}{N} \sum_i \log \left[\frac{p(t_i | \mathbf{z}_i, \theta^*)}{p(t_i | \mathbf{z}_i, \theta)} \right] \right\} \right\rangle_{\mathcal{D}} \\ &= - \left\langle \max_{\theta} \left\{ \frac{1}{N} \sum_i \log \left[\frac{p(t_i | \mathbf{z}_i, \theta)}{p(t_i | \mathbf{z}_i, \theta^*)} \right] \right\} \right\rangle_{\mathcal{D}} \\ &= - \lim_{\gamma \rightarrow \infty} \left\langle \frac{1}{\gamma} \log \int d\theta e^{\gamma \left\{ \frac{1}{N} \sum_i \log \left[\frac{p(t_i | \mathbf{z}_i, \theta)}{p(t_i | \mathbf{z}_i, \theta^*)} \right] \right\}} \right\rangle_{\mathcal{D}} \\ &= - \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma N} \left\langle \log \int d\theta e^{\gamma \sum_i \log \left[\frac{p(t_i | \mathbf{z}_i, \theta)}{p(t_i | \mathbf{z}_i, \theta^*)} \right]} \right\rangle_{\mathcal{D}} \quad \square \end{aligned}$$

interpretation:

stochastic minimisation, with noise $\sim 1/\gamma$

Step 3 – enable averaging over \mathcal{D}

▶ *Replica method*
(tame the log ...)

$$\langle \log Z \rangle = \lim_{n \rightarrow 0} \frac{1}{n} \log \langle Z^n \rangle = \lim_{n \rightarrow 0} \frac{1}{n} \log \left\langle \prod_{\alpha=1}^n Z \right\rangle$$

- evaluate for *integer* n ,
- analytical continuation to *non-integer* n

▶ *Application*

$$\begin{aligned} E(\theta^*) &= - \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma N} \left\langle \log \int d\theta e^{\gamma \sum_i \log [p(t_i | \mathbf{z}_i, \theta) / p(t_i | \mathbf{z}_i, \theta^*)]} \right\rangle_{\mathcal{D}} \\ &= - \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma N} \lim_{n \rightarrow 0} \frac{1}{n} \log \left\langle \left[\int d\theta e^{\gamma \sum_i \log [p(t_i | \mathbf{z}_i, \theta) / p(t_i | \mathbf{z}_i, \theta^*)]} \right]^n \right\rangle_{\mathcal{D}} \\ &= - \lim_{\gamma \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{\gamma N n} \log \int d\theta^1 \dots d\theta^n \left\langle e^{\gamma \sum_i \sum_{\alpha=1}^n \log [p(t_i | \mathbf{z}_i, \theta^\alpha) / p(t_i | \mathbf{z}_i, \theta^*)]} \right\rangle_{\mathcal{D}} \\ &= - \lim_{\gamma \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{\gamma N n} \log \int d\theta^1 \dots d\theta^n \left[\int d\mathbf{z} dt p(\mathbf{z}) p(t | \mathbf{z}, \theta^*) \prod_{\alpha=1}^n \left(\frac{p(t | \mathbf{z}, \theta^\alpha)}{p(t | \mathbf{z}, \theta^*)} \right)^\gamma \right]^N \end{aligned}$$

□

Track record of the replica method (Marc Kac, 1968)

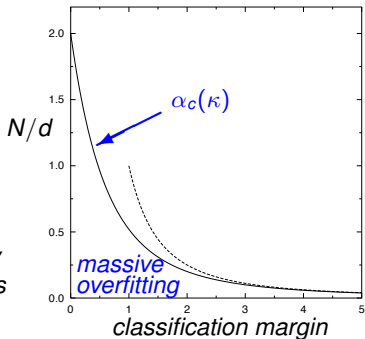
heterogeneous stochastic systems in physics,
biology, computer science, economics, ...

- ▶ *disordered magnets* (Sherrington & Kirkpatrick, 1975, Parisi, 1979)
- ▶ *attractor neural networks* (Amit, Gutfreund & Sompolinsky, 1985)
- ▶ *solution space of binary classifiers* (Gardner, 1988)

since then:

*satisfiability & optimisation problems,
error-correcting codes, minority games,
eigenvalue spectra of random graphs,
machine learning, protein folding,
immunology, compressed sensing, ...*

Gardner
theory
for binary
classifiers



Introduction

Regression for time-to-event data

Overfitting in PH regression

Replica analysis of overfitting in PH regression

The basic ideas

Translation to Cox's model

Replica symmetric solution

The end game

Simple approximations

Discussion

Translation to Cox's model

$$p(t|\mathbf{z}, \theta) \rightarrow p(t|\mathbf{z}, \lambda, \beta) = \lambda(t) e^{\beta \cdot \mathbf{z} / \sqrt{p} - \Lambda(t) \exp(\beta \cdot \mathbf{z} / \sqrt{p})}$$
$$\Lambda(t) = \int_0^t dt' \lambda(t')$$

- Defns, short-hands

$$p(\mathbf{z}) = (2\pi)^{-d/2} e^{-\frac{1}{2}\mathbf{z}^2}, \quad p(t|\xi, \lambda) = \lambda(t) e^{\xi - \Lambda(t) \exp(\xi)}$$
$$S^2 = \frac{1}{p} (\beta^*)^2, \quad \lambda^* = \lambda_0, \quad \alpha = \frac{d}{N}$$

- Insert, work out,
take $N \rightarrow \infty$:

$$E(S, \lambda^0) = - \lim_{\gamma \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{\gamma n} \text{extr}_{\{\mathbf{c}, \lambda_1, \dots, \lambda_n\}} \left\{ \frac{1}{2} \alpha n [1 + \log(2\pi)] + \frac{1}{2} \alpha \log \text{Det}(\mathbf{C}') \right.$$
$$\left. + \log \int \frac{d\mathbf{y} e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}}}{\sqrt{(2\pi)^{n+1} \text{Det} \mathbf{C}}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left(\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right)^\gamma \right\}$$

$$\mathbf{C}: \quad (n+1) \times (n+1), \quad C_{ab} = \langle \beta^a \cdot \beta^b / p \rangle, \quad a, b = 0 \dots n$$
$$\mathbf{C}': \quad n \times n, \quad C'_{ab} = C_{ab} - C_{a0} C_{0b} / C_{00}^2, \quad a, b = 1 \dots n$$

Introduction

Regression for time-to-event data

Overfitting in PH regression

Replica analysis of overfitting in PH regression

The basic ideas

Translation to Cox's model

Replica symmetric solution

The end game

Simple approximations

Discussion

Replica symmetric solution

If solution space connected:

saddle-point symmetric under *all* permutations of $\{1, \dots, n\}$

$$\mathbf{C} = \begin{pmatrix} S^2 & c_0 & \cdots & \cdots & c_0 \\ c_0 & C & c & \cdots & c \\ \vdots & c & C & \cdots & c \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_0 & c & \cdots & c & C \end{pmatrix}, \quad \lambda_\alpha(t) = \lambda(t) \quad \forall \alpha = 1 \dots n$$

interpretation:

$$C = \frac{1}{p} \langle \beta^2 \rangle_{\mathcal{D}}, \quad c = \frac{1}{p} \langle \beta \rangle_{\mathcal{D}}^2, \quad c_0 = \frac{1}{p} \beta^* \cdot \langle \beta \rangle_{\mathcal{D}}$$

Insert into formulae,
diagonalise \mathbf{C} and \mathbf{C}' , manipulations, integrations,
take the limit $n \rightarrow 0$...

Introduction

Regression for time-to-event data

Overfitting in PH regression

Replica analysis of overfitting in PH regression

The basic ideas

Translation to Cox's model

Replica symmetric solution

The end game

Simple approximations

Discussion

The end game ...

$$E_{RS}(S, \lambda_0) = - \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \text{extr}_{u, w; \lambda} \left\{ \alpha \log u \right. \\ \left. + \int \frac{dy_0}{\sqrt{2\pi}} e^{-\frac{1}{2}y_0^2} \int dt \rho(t|S y_0, \lambda_0) \log \int \frac{dy}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} \left[\frac{\rho(t|uy + wy_0, \lambda)}{\rho(t|S y_0, \lambda_0)} \right]^\gamma \right\}$$

□

interpretation:

$$u^2 = \frac{1}{\rho} \left(\langle \beta^2 \rangle_{\mathcal{D}} - \langle \beta \rangle_{\mathcal{D}}^2 \right), \quad w = \frac{1}{\rho} \langle \beta \rangle_{\mathcal{D}} \cdot \beta^* / |\beta^*|$$

- ▶ perfect fitting: $u \rightarrow 0$, $w \rightarrow S$, $\lambda(t) \rightarrow \lambda_0(t)$
indeed gives: $E_{RS}(S, \lambda_0) = 0$
provided $u \rightarrow 0$ slower than exponentially as $\gamma \rightarrow \infty$
- ▶ formula tricky, in view of
 - (i) non-commutation of limits $\alpha \rightarrow 0$ vs $u \rightarrow 0$ vs $\gamma \rightarrow \infty$
 - (ii) contains notorious integral

$$I(u, \kappa) = \int_{-\infty}^{\infty} dy e^{-\frac{1}{2}y^2 - \kappa \exp(uy)}$$

in Derrida 1981 for $u \rightarrow \infty$,
here we need $\kappa \rightarrow \infty \dots$

Simple approximations

- ▶ Approximate y -integral for finite $u > 0$

let $Dy = (2\pi)^{-1/2} e^{-\frac{1}{2}y^2} dy$:

$$E_{RS}(S, \lambda_0) \approx \lim_{\gamma \rightarrow \infty} \text{extr}_{\{u, w, \lambda\}} \left\{ \frac{1-\alpha}{\gamma} \log u + \frac{w^2}{2\gamma u^2} - \int Dy \int dt p(t|Sy, \lambda_0) \left[\log \left(\frac{\lambda(t)}{\lambda^*(t)\Lambda(t)} \right) - \frac{1}{\gamma u^2} \left(\frac{1}{2} \log^2 \Lambda(t) + wy \log \Lambda(t) \right) \right] \right\}$$

- ▶ Extremise over w and u

$$w = - \int Dy y \int dt p(t|Sy, \lambda_0) \log \Lambda(t)$$

$$u^2 = \frac{1}{1-\alpha} \left\{ \int Dy \int dt p(t|Sy, \lambda_0) \log^2 \Lambda(t) - \left[\int Dy y \int dt p(t|Sy, \lambda_0) \log \Lambda(t) \right]^2 \right\}$$

$$E_{RS}(S, \lambda_0) \approx \lim_{\gamma \rightarrow \infty} \text{extr}_{\{\lambda\}} \left\{ - \int Dy \int dt p(t|Sy, \lambda_0) \log \left(\frac{\lambda(t)}{\lambda^*(t)\Lambda(t)} \right) + \frac{1-\alpha}{2\gamma} \log \left[\int Dy \int dt p(t|Sy, \lambda_0) \log^2 \Lambda(t) - \left[\int Dy y \int dt p(t|Sy, \lambda_0) \log \Lambda(t) \right]^2 \right] \right\}$$

- ▶ Choose $\lambda_0(t) = \lambda_0$

Variational approx for inferred
base hazard rate: $\lambda(t) = \lambda$

results in

$$w = S, \quad u^2 = \frac{\pi^2/6}{1-\alpha}, \quad \lambda = \lambda_0 e^{C_E}$$

C_E : Euler's constant, ~ 0.5772

- ▶ everything scales sensibly with γ , limits can be taken
- ▶ u diverges for $\alpha \rightarrow 1$ (confirmed by simulations)
- ▶ $\lambda > \lambda_0$ (confirmed by simulations)

but approx of course too crude ...

- ▶ should have w and λ dependent on α ...
- ▶ should have $\lim_{\alpha \rightarrow 0} u = 0$...

Discussion

- ▶ Overfitting in PH causes predictable bias in regression pars
- ▶ Analytical approach to model overfitting in PH regression
- ▶ Problem reduced to analysis of an integral
- ▶ Preliminary approximations give rough but sensible results

- ▶ Next:
 - ▶ Tackle remaining integral properly
 - ▶ Generalise to correlated Gaussian covariates
 - ▶ Predict overfitting correction factors
 - ▶ Stochastic minimisation at optimal γ , such that $E = 0$

Thank you!