

Replica analysis of overfitting in survival analysis and other generalized linear models

ACC Coolen

King's College London and Saddle Point Science

Sendai, Oct 27th 2019

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Introduction

Data of modern medicine

Regression for time-to-event data

Overfitting in survival analysis

Phenomenology of overfitting

Failure of low-dimensional intuition

Quantitative theory of overfitting

Intuition for the problem

The basic ideas

The replica method

Applications of the theory

Cox regression

Regularized Cox regression

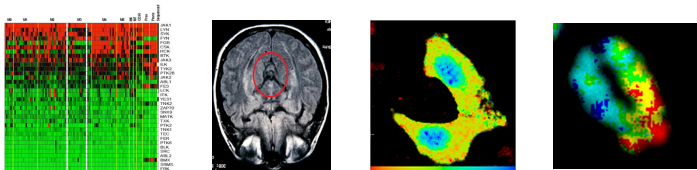
Extension to arbitrary generalized linear models

Generalized linear models

High-dimensional exponential inference models

Summary

Data of modern medicine



- ▶ *volume* ...
- ▶ *diversity* ...
(clinical, genomic, biomarkers, health records, imaging, ...)
- ▶ *complexity* of pipelines ...
(confounders, batch effects, variability between centres, ...)
- ▶ *dimensionality*
clinical ($\sim 10^1$), metabolic ($\sim 10^2$), gene expr ($\sim 10^5$),
imaging ($\sim 10^6$), NGS ($\sim 10^{10}$ and more)

generating 'big data' is not enough ...

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

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

\mathbf{z} : mutations, gene expr, biomarkers, images, ...

y : treatment response

- ▶ regression:

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

curse of dimensionality ...

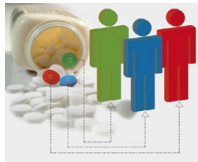
pre-genome: $N \sim 10^3$, $\dim \theta \sim 10^2$

post-genome: $N \sim 10^4$, $\dim \theta \sim 10^{10}$...

- ▶ simpler question: predict *individual* risk
(target aggressive treatments)

cancer: outcome is often a *duration* t ,
(overall or progression-free survival)

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



Introduction

Data of modern medicine

Regression for time-to-event data

Overfitting in survival analysis

Phenomenology of overfitting

Failure of low-dimensional intuition

Quantitative theory of overfitting

Intuition for the problem

The basic ideas

The replica method

Applications of the theory

Cox regression

Regularized Cox regression

Extension to arbitrary generalized linear models

Generalized linear models

High-dimensional exponential inference models

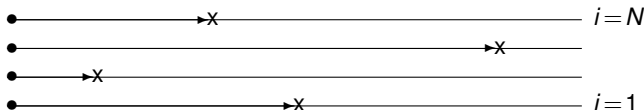
Summary

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)

\mathbf{z}_i : p covariates (measured at baseline)

$t_i > 0$: failure time (death, relapse, ...)



- ▶ *Aim*

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

1. identify high-risk patients
2. discover disease mechanisms
3. design new treatments (modifiable covariates)

Proportional hazards regression (1972)

main tool of medical statistics

$$\text{event time dist : } p(t|\mathbf{z}) = h(t|\mathbf{z})e^{-\int_0^t dt' h(t'|\mathbf{z})}$$

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

parameters: $\boldsymbol{\beta}, \lambda(t)$

- ▶ Maximum likelihood:

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

- ▶ Maximise over $\lambda(t)$ first

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

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

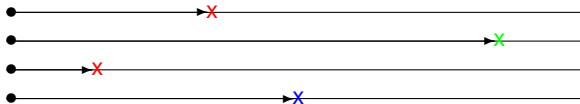
Beyond the basic model ...

► *Mathematical subtleties*

cure is forbidden, i.e. we need $\int_0^\infty dt \lambda(t) < \infty$,
violated by Breslow estimator ...

issues with event time coincidences,
i.e. we need $t_i \neq t_j$ for all i, j

► *Multiple competing risks*

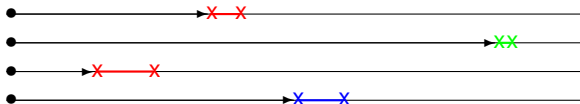


risk labels $r_i \in \{1, \dots, R\}$,

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

$$p(t, r | \mathbf{z}, \theta) = h_r(t | \mathbf{z}, \theta) \exp \left[- \int_0^t dt' \sum_{r'=1}^R h_{r'}(t' | \mathbf{z}, \theta) \right]$$

► *interval censoring*

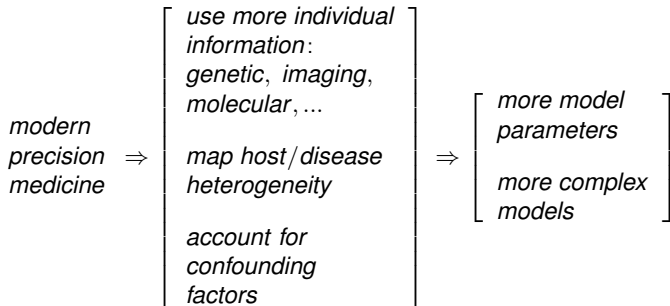


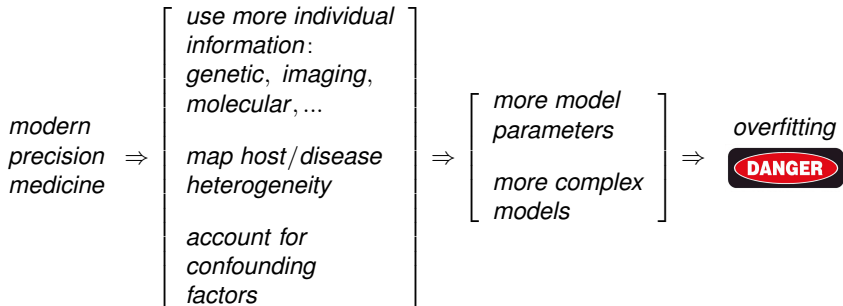
► *Include latent heterogeneity*

(e.g. mixture models, frailty models, random effects models)

distinct patient sub-classes, $\ell = 1 \dots L$,

$$p(t, r | \mathbf{z}, \boldsymbol{\theta}) = \sum_{\ell=1}^L w_{\ell} h_r(t' | \mathbf{z}, \boldsymbol{\theta}, \ell) \exp \left[- \int_0^t dt' \sum_{r'=1}^R h_{r'}(t' | \mathbf{z}, \boldsymbol{\theta}, \ell) \right]$$





Strategies to deal with overfitting in covariate-to-outcome analysis

- ▶ *'Back off'*

'safe' ratio covariates/samples
for multivariate regression

- ▶ *Eliminate redundant information*

improve covariates/samples ratio via
supervised nonlinear dimension reduction,
or by using biological knowledge.
'true' data dimension?

- ▶ *'Integrate out' overfitting effects*

fully Bayesian analysis of parameter uncertainty,
while keeping computation feasible

- ▶ *Model overfitting effects*

Overfitting correction theory for
multivariate regression



Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting**
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Phenomenology of overfitting

deteriorating outcome
prediction performance
on unseen data ...

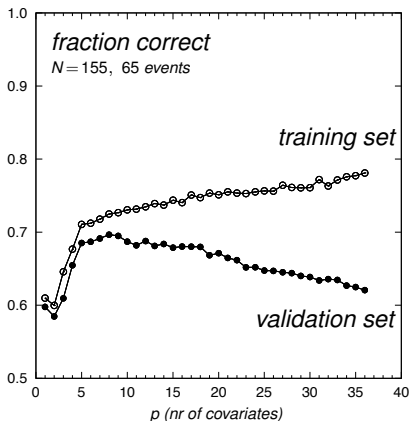
multivariate
Cox regression:

predict whether event before
or after a cutoff time point

primitive rule of thumb:

$$p_{\max} \sim \# \text{events} / 10$$

- ▶ too optimistic?
- ▶ indep of association strengths?
- ▶ indep of covariate correlations?



false positive associations ...

$N = 100$

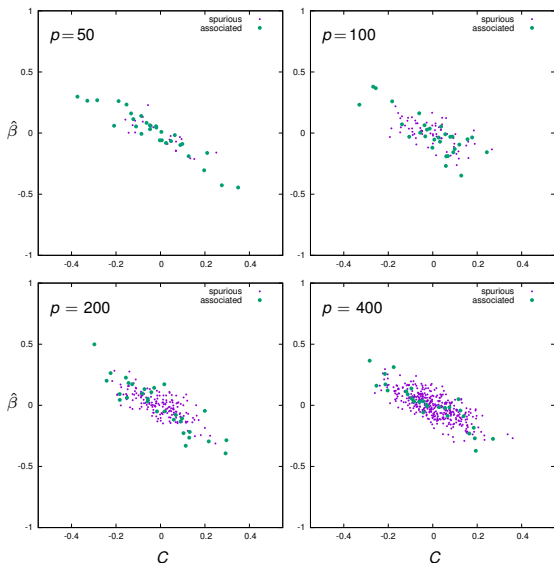
p covariates:

30 true associations

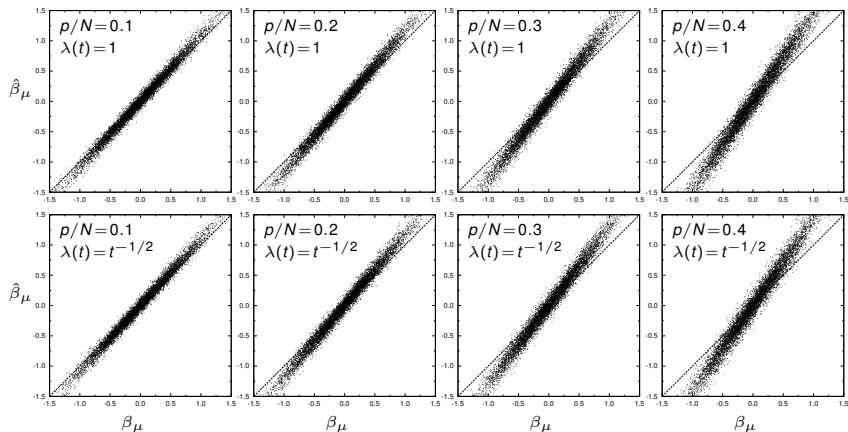
$p - 30$ spurious ones

C : Pearson correlation between covariates and event time,

$\hat{\beta}$: univariate Cox parameters



bias in inferred association parameters ...



β_μ : true associations

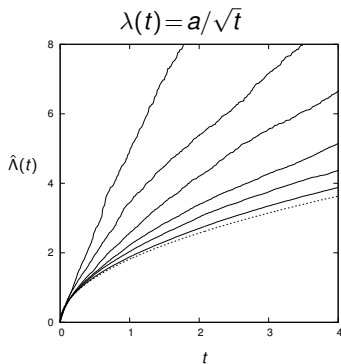
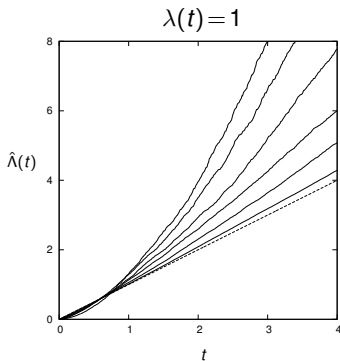
$\hat{\beta}_\mu$: multivariate regression

synthetic survival data, generated from Cox model with $N=400$

bias in inferred base hazard rates ...

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

$$p/N = 0.05 \rightarrow 0.55$$



synthetic survival data,
generated from Cox model with $N = 400$

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition**

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Failure of our low-dimensional intuition

low-dim ML/MAP regime: $N \rightarrow \infty$, p fixed

high-dim regime: $N, p \rightarrow \infty$, p/N finite

- ▶ hyperparameters of priors (regularizers) must be p -dependent ...
- ▶ for sufficiently large N :
regression possible even when $p > N$...

notation:

data : $\mathcal{D} = \{(\mathbf{z}_1, t_1), \dots, (\mathbf{z}_N, t_N)\}$, $\mathbf{z}_i \in \mathbb{R}^p$ *covariates*
 $t_i > 0$ *event time*

Cox model : $p(t|\mathbf{z}, \beta, \lambda) = -\frac{d}{dt} \exp[-e^{\beta \cdot \mathbf{z}} \Lambda(t)]$

ML inference : $(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \left\{ \sum_{i=1}^N \log p(t_i | \mathbf{z}_i, \beta, \lambda) \right\}$

MAP inference : $(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \left\{ \sum_{i=1}^N \log p(t_i | \mathbf{z}_i, \beta, \lambda) + \log \overbrace{p(\beta)}^{\text{prior}} \right\}$

scaling of hyperparameters

in regularised Cox regression

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

$$p(t | \mathbf{z}, \beta, \lambda) = -\frac{d}{dt} \exp[-e^{\beta \cdot \mathbf{z}} \Lambda(t)], \quad \beta \cdot \mathbf{z} = \sum_{\mu=1}^p \beta_{\mu} z_{\mu}$$

e.g. $p(\beta) \propto e^{-\sum_{\mu=1}^p |\beta_{\mu}/\sigma|}$, $p(\beta) \propto e^{-\frac{1}{2} \sum_{\mu=1}^p (\beta_{\mu}/\sigma)^2}$, σ : hyperpar

claim: $\sigma = \mathcal{O}(p^{-\frac{1}{2}})$ as $p \rightarrow \infty$

- ▶ general theory
- ▶ simple scaling argument:

$$\beta_{\mu} = \mathcal{O}(1): \quad \sum_{\mu=1}^p \beta_{\mu} z_{\mu}^i = \mathcal{O}(\sqrt{p}) \Rightarrow t_i \in \{0, \infty\}$$

$$\beta_{\mu} = \mathcal{O}(p^{-1/2}): \quad \sum_{\mu=1}^p \beta_{\mu} z_{\mu}^i = \mathcal{O}(1) \Rightarrow t_i \text{ finite}$$

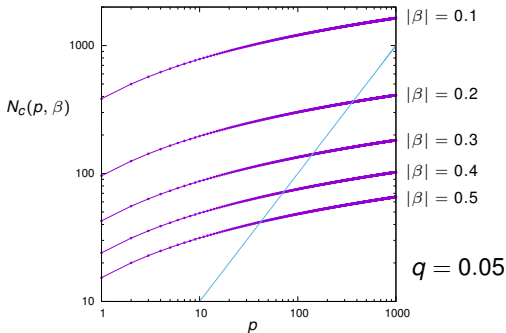
finite event times = prior knowledge!

regression with $p > N$ in principle possible
provided N is large enough

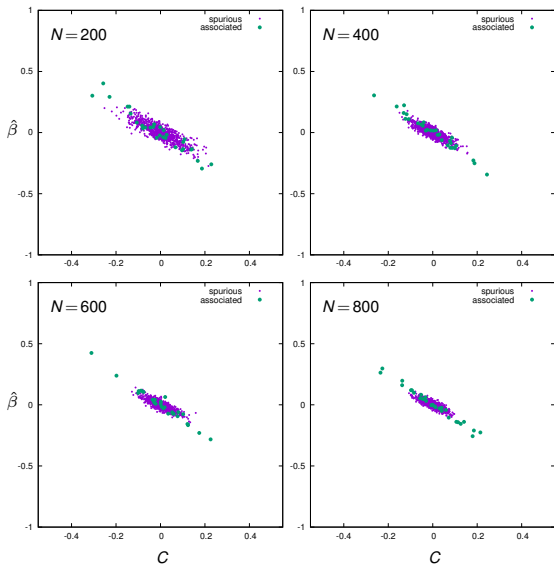
- ▶ $N \uparrow$: prob of false positive associations \downarrow
- ▶ $p \uparrow$: prob of false positive associations \uparrow

uncorr covars: $N > N_c(p, \beta)$: prob of finding one or more spurious univariate associations of strength $\geq |\beta|$ is less than q

$$N_c(p, \beta) = \frac{2}{\beta^2} \left[\text{Erf}^{-1} \left(e^{\frac{1}{p} \log(1-q)} \right) \right]^2$$



synthetic survival data,
with $p = 800$



Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

Intuition for the problem

- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Intuition for the problem

- ▶ information-theoretic interpretation of ML regression

assumed model: p_{θ}

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

$$\hat{p}(t, \mathbf{z}) = \frac{1}{N} \sum_{i=1}^N \delta(t-t_i) \delta(\mathbf{z}-\mathbf{z}_i) \quad (\text{empirical distribution})$$

$$D(\hat{p}||p_{\theta}) = \int dt d\mathbf{z} \hat{p}(t, \mathbf{z}) \log \left[\frac{\hat{p}(t|\mathbf{z})}{p(t|\mathbf{z}, \theta)} \right] \quad (\text{KL-distance})$$

- ▶ so ML regression pushes $p(t|\mathbf{z}, \theta)$ as close as possible towards $\hat{p}(t|\mathbf{z})$

true pars: θ^*

- ▶ fixed p , $N \rightarrow \infty$: $\hat{p}(t, \mathbf{z}) = p(t, \mathbf{z}|\theta^*)$, so $\theta_{\text{ML}} = \theta^*$ ✓
- ▶ $p = \mathcal{O}(N)$, $N \rightarrow \infty$: $\hat{p}(t, \mathbf{z}) \neq p(t, \mathbf{z}|\theta^*)$, so $\theta_{\text{ML}} \neq \theta^*$ ✗

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem

The basic ideas

- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

The basic ideas

Step1 – identify quantity to calculate

- ▶ \hat{p}_{θ^*} : empirical distr of (t, \mathbf{z}) ,
for data generated with θ^*

ML regression: minimize $D(\hat{p}_{\theta^*} || p_{\theta})$

optimal stopping point: $\theta = \theta^*$

$$D(\hat{p}_{\theta^*} || p_{\theta}) = D(\hat{p}_{\theta^*} || p_{\theta^*}) \quad \leftarrow \text{zero iff } p \ll N$$

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

$E(\theta^*, \mathcal{D}) < 0$: 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 – remove minimisation over θ

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

► *Laplace identity*

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

use in reverse:

$$E(\theta^*) = - \lim_{\gamma \rightarrow \infty} \frac{1}{N} \frac{\partial}{\partial \gamma} \left\langle \log \int d\theta \prod_{i=1}^N \left[\frac{\rho(t_i | \mathbf{z}_i, \theta)}{\rho(t_i | \mathbf{z}_i, \theta^*)} \right]^{\gamma} \right\rangle_{\mathcal{D}}$$

interpretation:

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

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method**

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

The replica method

aim: make hard analytical calculations easy ...

here: compute the average over \mathcal{D}

$$E(\theta^*) = - \lim_{\gamma \rightarrow \infty} \frac{1}{N} \frac{\partial}{\partial \gamma} \left\langle \log \int d\theta \prod_{i=1}^N \left[\frac{\rho(t_i | \mathbf{z}_i, \theta)}{\rho(t_i | \mathbf{z}_i, \theta^*)} \right]^\gamma \right\rangle_{\mathcal{D}}$$

- ▶ replica method

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

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

- ▶ application

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

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression**
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Translation to Cox's model

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

- ▶ since $\langle z_\mu \rangle$ and $\langle z_\mu z_\nu \rangle$ can be transformed away, definitions:

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

$$S^2 = (\beta^*)^2 / p, \quad \lambda^* = \lambda_0, \quad \alpha = d/N$$

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

$$E(S, \lambda_0) = - \lim_{\gamma \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{n} \frac{\partial}{\partial \gamma} \text{extr}_{\{\mathbf{C}, \lambda_1, \dots, \lambda_n\}} \left\{ \frac{1}{2} \alpha n [1 + \log(2\pi)] + \frac{1}{2} \alpha \log \text{Det}(\mathbf{C}') \right.$$

$$\left. + \log \int \frac{d\mathbf{y}}{\sqrt{(2\pi)^{n+1} \text{Det} \mathbf{C}}} e^{-\frac{1}{2} \mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left(\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right)^\gamma \right\}$$

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

$$\mathbf{C}': \quad n \times n, \quad C'_{ab} = C_{ab} - C_{a0} C_{0b} / C_{00}^2, \quad a, b = 1 \dots n$$

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:

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

Insert into formulae,
diagonalise \mathbf{C} and \mathbf{C}' ,
manipulations, integrations,
take limits $n \rightarrow 0$ and $\gamma \rightarrow \infty \dots$

- ▶ explicit formula for $E(S, \lambda^*, \zeta)$

$$\zeta = p/N, \quad S = |\beta^*|$$

$\beta^*, \lambda^*(t)$: true associations and base hazard rate

- ▶ requires solving $u, v, w, \lambda(t)$ from

$$\zeta v^2 = \int DzDy \int dt p(t|Sy, \lambda^*) \left[u^2 - W(u^2 e^{u^2 + wy + vz} \Lambda(t)) \right]^2$$

$$\zeta = \int DzDy \int dt p(t|Sy, \lambda^*) \frac{W(u^2 e^{u^2 + wy + vz} \Lambda(t))}{1 + W(u^2 e^{u^2 + wy + vz} \Lambda(t))}$$

$$0 = \int DzDy y \int dt p(t|Sy, \lambda^*) W(u^2 e^{u^2 + wy + vz} \Lambda(t))$$

$$\frac{p(t)}{\lambda(t)} = \int DzDy \int_t^\infty dt' p(t'|Sy, \lambda^*) \frac{W(u^2 e^{u^2 + wy + vz} \Lambda(t'))}{u^2 \Lambda(t')}$$

$$Dz = (2\pi)^{-1/2} e^{-\frac{1}{2}z^2} dz$$

$W(x)$: Lambert function

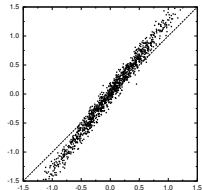
$$p(t|\xi, \lambda) = \lambda(t) e^{\xi - \exp(\xi)\Lambda(t)}$$

- ▶ interpretation:

$$\text{slope} : \kappa = w/S$$

$$\text{width} : \sigma = v/\sqrt{\rho}$$

*all we need for
overfitting correction!*



- ▶ challenge: eqn for $\Lambda(t)$

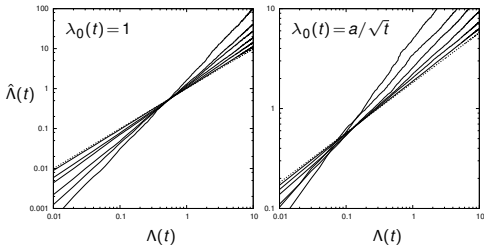
$$t \gg 1 : \quad \log \Lambda(t) = \rho \log \Lambda^*(t) + (1-\rho) \log \log \Lambda^*(t) + \dots$$

$$\rho = \frac{w}{2S} \left(1 + \sqrt{1 + 4u^2/w^2} \right)$$

- ▶ variational approx:

$$\Lambda(t) = k[\Lambda^*(t)]^\rho$$

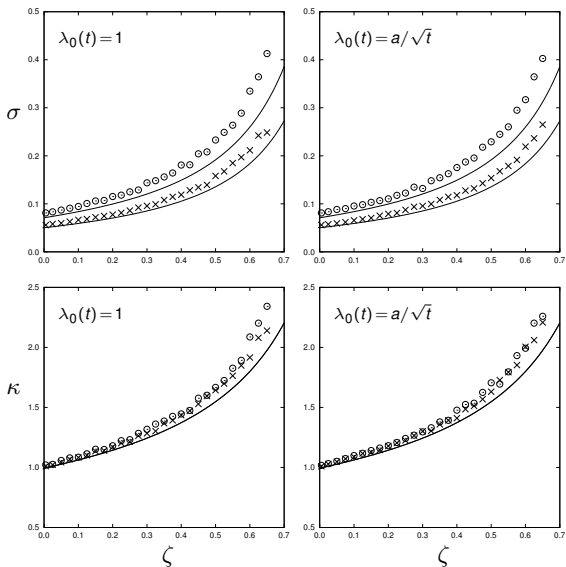
$\Lambda^*(t)$ drops out of equations!



width σ and slope κ
of data clouds

lines: variational theory
for $S = 0.5$ and $\langle t \rangle = 1$

simulations:
 \circ : $N = 200$
 \times : $N = 400$



Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression**

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Application to regularized Cox regression

$$N, p \rightarrow \infty, \quad \zeta = p/N,$$

$$\text{L2-prior: } \rho(\beta) \propto \exp(-\eta p \beta^2)$$

main changes:

- ▶ dependence of theory on eigenvalue spectrum $\varrho(a)$ of covariate correlation matrix
- ▶ two extra order parameters, closed equations for $u, v, w, f, g, \lambda(t)$
- ▶ no longer a phase transition at $\zeta = 1$, but inference well-defined for any $\zeta > 0$
- ▶ closed eqn for *optimal* hyperparameter η , defined by demanding absent bias, i.e. slope $\kappa = 1$

scalar order parameter eqns:

$$\zeta \tilde{f} \tilde{u}^4 = - \int DzDy \int dt \rho(t|S\langle a \rangle^{\frac{1}{2}} y, \lambda^*) \left[W(\tilde{u}^2 e^{\tilde{u}^2 + wy + vz} \Lambda(t)) - \tilde{u}^2 \right]^2$$

$$\zeta \tilde{g} \tilde{u}^2 = \int DzDy \int dt \rho(t|S\langle a \rangle^{\frac{1}{2}} y, \lambda^*) \frac{W(\tilde{u}^2 e^{\tilde{u}^2 + wy + vz} \Lambda(t))}{1 + W(\tilde{u}^2 e^{\tilde{u}^2 + wy + vz} \Lambda(t))}$$

$$0 = \zeta w \left[\langle a \rangle \left\langle \frac{a^2}{2\eta + \tilde{g}a} \right\rangle^{-1} - \tilde{g} \right] + \frac{1}{\tilde{u}^2} \int DzDy y \int dt \rho(t|S\langle a \rangle^{\frac{1}{2}} y, \lambda^*) W(\tilde{u}^2 e^{\tilde{u}^2 + wy + vz} \Lambda(t))$$

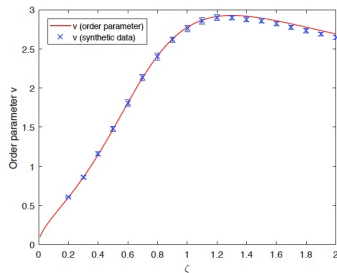
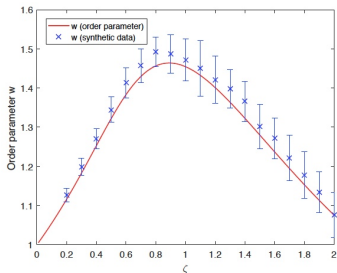
$$\tilde{u}^2 = \left\langle \frac{a}{2\eta + \tilde{g}a} \right\rangle$$

$$v^2 = w^2 \left[\langle a \rangle \left\langle \frac{a^2}{2\eta + \tilde{g}a} \right\rangle^{-2} \left\langle \frac{a^3}{(2\eta + \tilde{g}a)^2} \right\rangle - 1 \right] - \tilde{f} \left\langle \frac{a^2}{(2\eta + \tilde{g}a)^2} \right\rangle$$

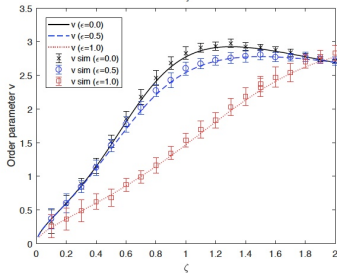
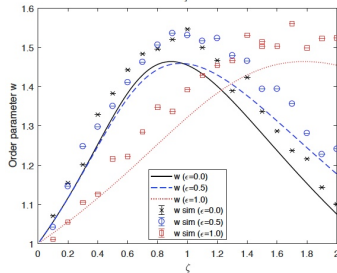
functional order parameter eqn:

$$\frac{\rho(t)}{\lambda(t)} = \int DzDy \int_t^\infty \frac{dt'}{\tilde{u}^2 \Lambda(t')} \rho(t'|S\langle a \rangle^{\frac{1}{2}} y, \lambda^*) W(\tilde{u}^2 e^{\tilde{u}^2 + wy + vz} \Lambda(t'))$$

*uncorrelated
covariates*



*pairwise
correlated
covariates*



$\eta = 0.025$, top row $p=2000$, bottom row $Np=400,000$

overfitting correction pars: slope $\kappa = w/\tilde{S}$, width $\sigma = v/\sqrt{\bar{p}}$

final overfitting correction protocol

1. estimate covariate correlation matrix \mathbf{A} ,
compute its eigenvalue spectrum $\varrho(a)$
2. carry out regularized Cox regression,
with prior $p(\beta) \propto \exp(-\eta p\beta^2)$ (small $\eta > 0$)
result: $\hat{\beta}$ and $\hat{\Lambda}(t)$ (Breslow estimator)
3. calculate $\hat{\beta} \cdot \mathbf{A}\hat{\beta}$
4. solve coupled nonlinear eqns for (u, v, w, f, g, k, ρ) ,
alongside $v^2 + w^2 = \hat{\beta} \cdot \mathbf{A}\hat{\beta}$
(replaces unknown variable S)
5. calculate slope κ and noise amplitude σ
6. compute corrected estimators:

$$\hat{\beta} = \kappa^{-1} \hat{\beta}, \quad \hat{\Lambda}(t) = [\hat{\Lambda}(t)/k]^{1/\rho}$$

7. use σ to correct p -values

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models**
- High-dimensional exponential inference models

Summary

Generalized linear models

$$p(y|\mathbf{z}, \theta) = p(y|\beta^1 \cdot \mathbf{z}, \dots, \beta^K \cdot \mathbf{z}; \omega)$$

- ▶ *Logistic regression*

$$p(y|\mathbf{z}, \{\beta\}, \beta_0) = \frac{e^{y(\beta \cdot \mathbf{z} + \beta_0)}}{2 \cosh(\beta \cdot \mathbf{z} + \beta_0)}, \quad y = \pm 1$$

- ▶ *Ordinal class regression*

$$p(y|\mathbf{z}, \{\beta\}, \lambda) = e^{-\exp(\beta \cdot \mathbf{z}) \sum_{y' > y} \lambda_{y'}} - e^{-\exp(\beta \cdot \mathbf{z}) \sum_{y' \geq y} \lambda_{y'}}, \quad y = 1, \dots, C$$

- ▶ *Latent class survival analysis*

$$p(t|\mathbf{z}, \{\beta\}, \{\lambda\}, \mathbf{w}) = \sum_{\ell=1}^K w_{\ell} \left[\lambda_{\ell}(t) e^{\beta^{\ell} \cdot \mathbf{z} - \exp(\beta^{\ell} \cdot \mathbf{z}) \int_0^t dt' \lambda_{\ell}(t')} \right], \quad t > 0$$

- ▶ *Neural networks*

$$p(y|\mathbf{z}, \{\beta\}, \mathbf{w}, \mathbf{v}) = \operatorname{sgn} \left[\sum_{\ell=1}^K w_{\ell} \operatorname{sgn} \left[\beta^{\ell} \cdot \mathbf{z} + v_{\ell} \right] + v_0 \right], \quad y = \pm 1$$

- ▶

$$K = 1 : \quad p(s|\mathbf{z}, \theta) = p(s|\beta \cdot \mathbf{z}, \theta), \quad p(\beta) \propto \exp(-\eta p \beta^2)$$

replica analysis (RS):

order pars $(\tilde{f}, \tilde{g}, \tilde{u}, v, w, \theta)$

$$\left\langle \frac{a}{2\eta + \tilde{g}a} \right\rangle = \tilde{u}^2$$

$$w^2 \left[\langle a \rangle \left\langle \frac{a^2}{2\eta + \tilde{g}a} \right\rangle^{-2} \left\langle \frac{a^3}{(2\eta + \tilde{g}a)^2} \right\rangle - 1 \right] - \tilde{f} \left\langle \frac{a^2}{(2\eta + \tilde{g}a)^2} \right\rangle = v^2$$

$$\int \text{DyDz} \langle [\xi(\mathbf{w}\mathbf{y} + \mathbf{v}\mathbf{z}, \tilde{u}, \mathbf{s}, \theta) - \mathbf{w}\mathbf{y} - \mathbf{v}\mathbf{z}]^2 \rangle_s = -\zeta \tilde{f} \tilde{u}^4$$

$$\int \text{DyDz} \langle (\partial_1 \xi)(\mathbf{w}\mathbf{y} + \mathbf{v}\mathbf{z}, \tilde{u}, \mathbf{s}, \theta) \rangle_s = 1 - \zeta \tilde{g} \tilde{u}^2$$

$$\int \text{DyDz} \left\langle \xi(\mathbf{w}\mathbf{y} + \mathbf{v}\mathbf{z}, \tilde{u}, \mathbf{s}, \theta) \frac{\partial \log p(\mathbf{s} | \mathcal{S}\langle a \rangle^{\frac{1}{2}} \mathbf{y}, \theta^*)}{\partial \mathbf{y}} \right\rangle_s = \zeta w \tilde{u}^2 \langle a \rangle \left\langle \frac{a^2}{2\eta + \tilde{g}a} \right\rangle^{-1}$$

$$\int \text{DyDz} \left\langle \frac{\partial \log p(\mathbf{s} | \xi(\mathbf{w}\mathbf{y} + \mathbf{v}\mathbf{z}, \tilde{u}, \mathbf{s}, \theta), \theta)}{\partial \theta} \right\rangle_s = 0$$

with

$$\langle F(\mathbf{s}) \rangle_s = \int \text{d}\mathbf{s} p(\mathbf{s} | \mathcal{S}\langle a \rangle^{\frac{1}{2}} \mathbf{y}, \theta^*) F(\mathbf{s}) \quad \text{or} \quad \sum_s p(\mathbf{s} | \mathcal{S}\langle a \rangle^{\frac{1}{2}} \mathbf{y}, \theta^*) F(\mathbf{s})$$

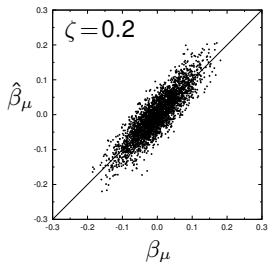
$$\xi(\mu, \sigma, \mathbf{s}, \theta) = \operatorname{argmax}_{\xi \in \mathbb{R}} \left[\log p(\mathbf{s} | \xi, \theta) - \frac{1}{2} (\xi - \mu)^2 / \sigma^2 \right]$$

Test 1

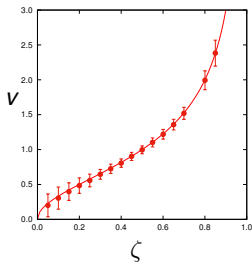
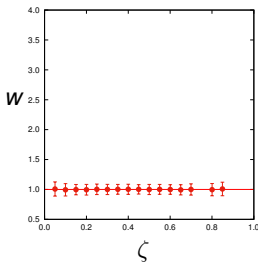
- ▶ *Linear regression* $s \in \mathbb{R}$, $p(s|\beta \cdot \mathbf{z}, \beta_0) = e^{-\frac{1}{2}(s - \beta \cdot \mathbf{z} - \beta_0)^2 / \Sigma^2} / \Sigma \sqrt{2\pi}$

prediction for ML ($\eta = 0$):

$$\begin{aligned} \beta_0 &= \beta_0^*, & w &= S\langle a \rangle^{\frac{1}{2}}, & \Sigma &= \Sigma^* \sqrt{1 - \zeta} \quad \checkmark \\ \tilde{u} &= \Sigma^* \sqrt{\zeta}, & v &= \Sigma^* \left(\frac{\zeta}{1 - \zeta} \right)^{\frac{1}{2}} \end{aligned}$$



$Np = 400,000$



no overfitting-induced bias,
but under-estimation of model noise

Test 2

▶ *Cox model*

$$t \geq 0, \quad p(t|\beta \cdot \mathbf{z}, \lambda) = \lambda(t) e^{\beta \cdot \mathbf{z} - \exp(\beta \cdot \mathbf{z}) \int_0^t \lambda(t') dt'}$$

recover fully all previous

MAP/ML overfitting equations



Application
to logistic
regression

$$s = \pm 1, \quad p(s|\beta \cdot \mathbf{z}, \beta_0) = \frac{e^{s(\beta \cdot \mathbf{z} + \beta_0)}}{2 \cosh(\beta \cdot \mathbf{z} + \beta_0)}$$

$$\zeta v^2 = \int DyDz \left\{ \frac{1}{2} \left[1 + \tanh(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \left[\tilde{x}(\beta_0 + \mathbf{w}y + \mathbf{v}z, \tilde{u}) - (\beta_0 + \mathbf{w}y + \mathbf{v}z) \right]^2 \right. \\ \left. + \frac{1}{2} \left[1 - \tanh(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \left[\tilde{x}(-(\beta_0 + \mathbf{w}y + \mathbf{v}z), \tilde{u}) + (\beta_0 + \mathbf{w}y + \mathbf{v}z) \right]^2 \right\}$$

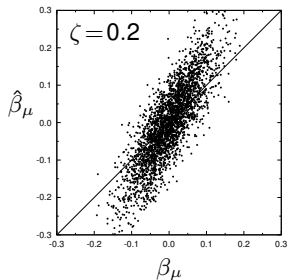
$$\zeta = \int DyDz \left\{ \frac{1}{2} \left[1 + \tanh(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \frac{\tilde{u}^2 [1 - \tanh^2(\tilde{x}(\beta_0 + \mathbf{w}y + \mathbf{v}z, \tilde{u}))]}{1 + \tilde{u}^2 [1 - \tanh^2(\tilde{x}(\beta_0 + \mathbf{w}y + \mathbf{v}z, \tilde{u}))]} \right. \\ \left. + \frac{1}{2} \left[1 - \tanh(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \frac{\tilde{u}^2 [1 - \tanh^2(\tilde{x}(-(\beta_0 + \mathbf{w}y + \mathbf{v}z), \tilde{u}))]}{1 + \tilde{u}^2 [1 - \tanh^2(\tilde{x}(-(\beta_0 + \mathbf{w}y + \mathbf{v}z), \tilde{u}))]} \right\}$$

$$\zeta \mathbf{w} = \frac{1}{2} S\langle a \rangle^{\frac{1}{2}} \int DyDz \left[1 - \tanh^2(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \\ \times \left[\tilde{x}(\beta_0 + \mathbf{w}y + \mathbf{v}z, \tilde{u}) + \tilde{x}(-(\beta_0 + \mathbf{w}y + \mathbf{v}z), \tilde{u}) \right]$$

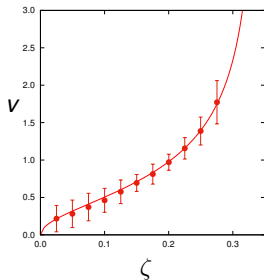
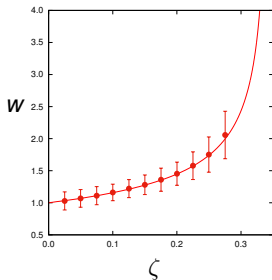
$$\beta_0 = \int DyDz \left\{ \frac{1}{2} \left[1 + \tanh(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \tilde{x}(\beta_0 + \mathbf{w}y + \mathbf{v}z, \tilde{u}) \right. \\ \left. - \frac{1}{2} \left[1 - \tanh(S\langle a \rangle^{\frac{1}{2}} y + \beta_0^*) \right] \tilde{x}(-(\beta_0 + \mathbf{w}y + \mathbf{v}z), \tilde{u}) \right\}$$

with

$$\tilde{X}(\eta, \sigma) : \quad \text{soln of } \tanh(x) = 1 - (x - \eta)/\sigma^2$$



$Np = 400,000$



w : slope of $\{\hat{\beta}_\mu\}$ cloud
 v/\sqrt{p} : width of $\{\hat{\beta}_\mu\}$ cloud

Introduction

- Data of modern medicine
- Regression for time-to-event data

Overfitting in survival analysis

- Phenomenology of overfitting
- Failure of low-dimensional intuition

Quantitative theory of overfitting

- Intuition for the problem
- The basic ideas
- The replica method

Applications of the theory

- Cox regression
- Regularized Cox regression

Extension to arbitrary generalized linear models

- Generalized linear models
- High-dimensional exponential inference models

Summary

Overfitting in high-dimensional exponential inference models

observed data: $\mathbf{x}_i \in A^p$,

assumed model: $p(\mathbf{x}|\boldsymbol{\theta}^*)$, $\boldsymbol{\theta}^* \in \mathbb{R}^q$

- ▶ Exponential models:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{e^{\boldsymbol{\theta} \cdot \boldsymbol{\omega}(\mathbf{x})/\sqrt{q}}}{|A|^p Z(\boldsymbol{\theta})}, \quad Z(\boldsymbol{\theta}) = \frac{1}{|A|^p} \sum_{\mathbf{x}} e^{\boldsymbol{\theta} \cdot \boldsymbol{\omega}(\mathbf{x})/\sqrt{q}}$$

e.g: inference of

- spin interactions from observed configurations
- aminoacid contact maps from protein structures

- ▶ MAP inference, with Gaussian priors:

$$E_\gamma(\boldsymbol{\theta}^*) = -\frac{\boldsymbol{\theta}^{*2}}{2N\sigma^2} - \frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{Nn} \log \int d\boldsymbol{\theta}^1 \dots d\boldsymbol{\theta}^n \left[\prod_{\alpha=1}^n e^{-\frac{1}{2}\gamma \boldsymbol{\theta}^{\alpha 2}/\sigma^2} \right] e^{N\psi(\{\boldsymbol{\theta}^\alpha\}, \boldsymbol{\theta}^*)}$$

$$\Psi(\{\boldsymbol{\theta}^\alpha\}, \boldsymbol{\theta}^*) = \log Z\left(\boldsymbol{\theta}^* + \gamma \sum_{\alpha=1}^n (\boldsymbol{\theta}^\alpha - \boldsymbol{\theta}^*)\right) - \gamma \sum_{\alpha=1}^n \log Z(\boldsymbol{\theta}^\alpha) - (1 - \gamma n) \log Z(\boldsymbol{\theta}^*)$$

Summary

- ▶ *Overfitting in Cox regression*

bias in regression parameters: $\hat{\beta} = \kappa\beta^* + \text{noise}$, $\kappa > 1$
bias in base hazard rates

- ▶ *Analytical approach based on the replica method*

RS: closed equations for $\{u, v, w, \lambda(t)\}$
variational approximation for $\lambda(t)$
correct predictions for slope κ and noise σ

- ▶ *L2-regularized Cox regression*

theory involves spectrum of covariate correlation matrix
 p -dependent hyperparameter
phase transition removed

- ▶ *Generalizations of the theory*

arbitrary generalized linear regression models
arbitrary high-dimensional exponential models

Papers, talks, seminars

<https://nms.kcl.ac.uk/ton.coolen>

Coolen et al, J Phys. A 50 (2017)

Sheikh and Coolen, J. Phys. A 52 (2019)

Thanks to

collaborators: James Barrett, Alexander Mozeika, Pierre Paga,
Conrad Perez-Vicente, Mansoor Sheikh

funding: MRC, CRUK, GSK

January 2020



King's College London



Radboud University

