

Overfitting correction in multivariate survival analysis

ACC Coolen

King's College London and Saddle Point Science

Belfast, Sept 4th 2019

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

- Other extensions

Summary

*modern
precision
medicine*

⇒

*use more individual
information:
genetic, imaging,
molecular, ...*

*map host/disease
heterogeneity*

*account for
confounding
factors*

⇒

*more model
parameters*

*more complex
models*

*modern
precision
medicine*

⇒

*use more individual
information:
genetic, imaging,
molecular, ...*

*map host/disease
heterogeneity*

*account for
confounding
factors*

⇒

*more model
parameters*

*more complex
models*

*modern
precision
medicine*

⇒

*use more individual
information:
genetic, imaging,
molecular, ...*

*map host/disease
heterogeneity*

*account for
confounding
factors*

⇒

*more model
parameters*

*more complex
models*

*modern
precision
medicine*



*use more individual
information:
genetic, imaging,
molecular, ...*

*map host/disease
heterogeneity*

*account for
confounding
factors*



*more model
parameters*

*more complex
models*



overfitting

A red oval with the word "DANGER" in white capital letters on a black background.

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

Other extensions

Summary

Phenomenology of overfitting

deteriorating outcome
prediction performance
on unseen data ...

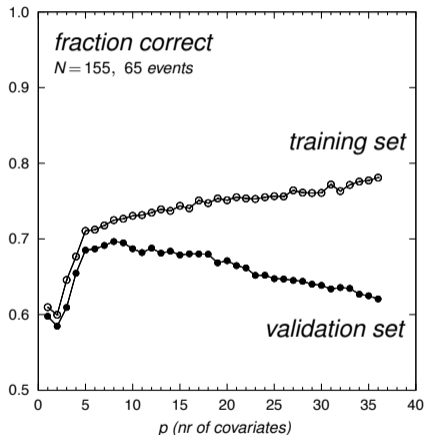
multivariate
Cox regression:

predict whether event before
or after a cutoff time point

primitive rule of thumb:

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

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



Phenomenology of overfitting

deteriorating outcome
prediction performance
on unseen data ...

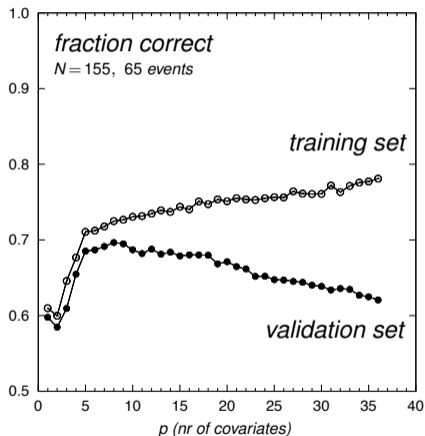
multivariate
Cox regression:

predict whether event before
or after a cutoff time point

primitive rule of thumb:

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

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



false positive associations ...

$N = 100$

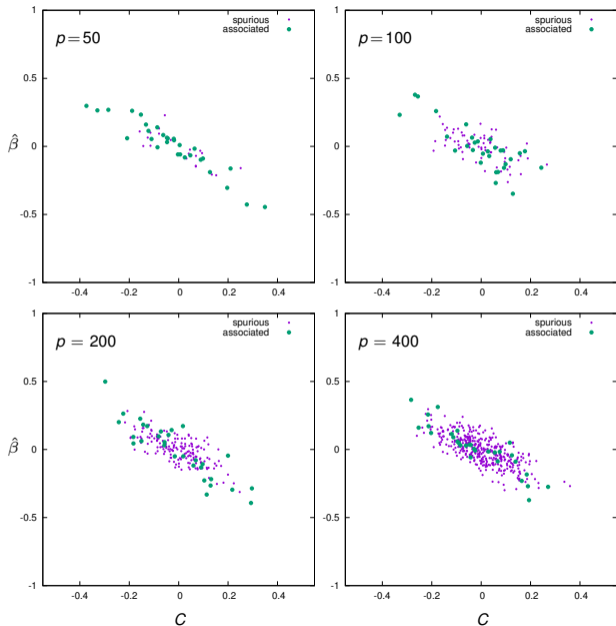
p covariates:

30 true associations

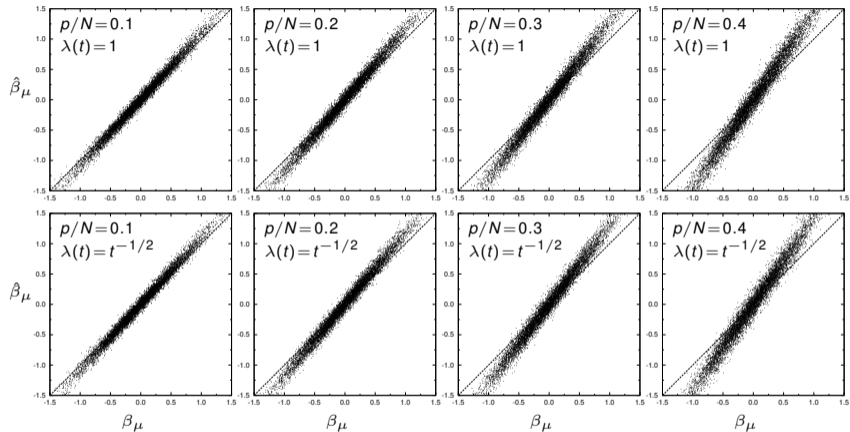
$p - 30$ spurious ones

C : Pearson correlation between covariates and event time,

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



bias in inferred association parameters ...



β_μ : true associations

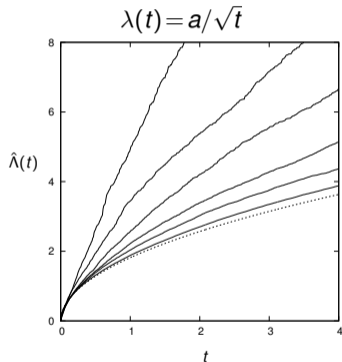
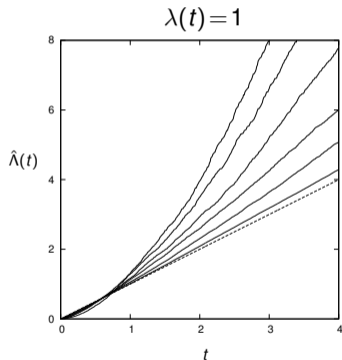
$\hat{\beta}_\mu$: multivariate regression

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

bias in inferred base hazard rates ...

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

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



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

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

Other extensions

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

- ▶ p -values are blind to overfitting ...
- ▶ 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 : $(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \left\{ \sum_{i=1}^N \log p(t_i | \mathbf{z}_i, \beta, \lambda) + \log \overbrace{p(\beta)}^{\text{prior}} \right\}$

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

- ▶ p -values are blind to overfitting ...
- ▶ 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!

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!

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:

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

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

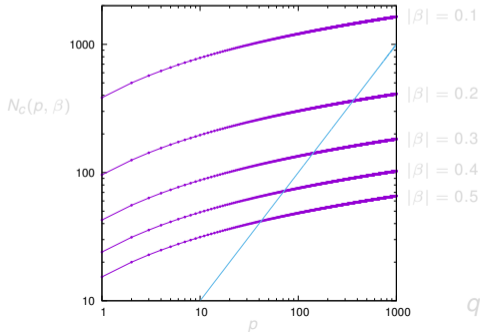
finite event times = prior knowledge!

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

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

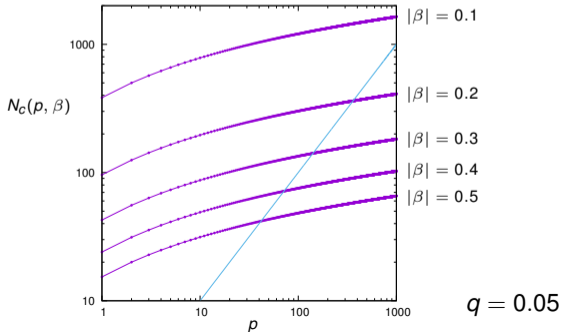


regression with $p > N$ in principle possible
provided N is large enough

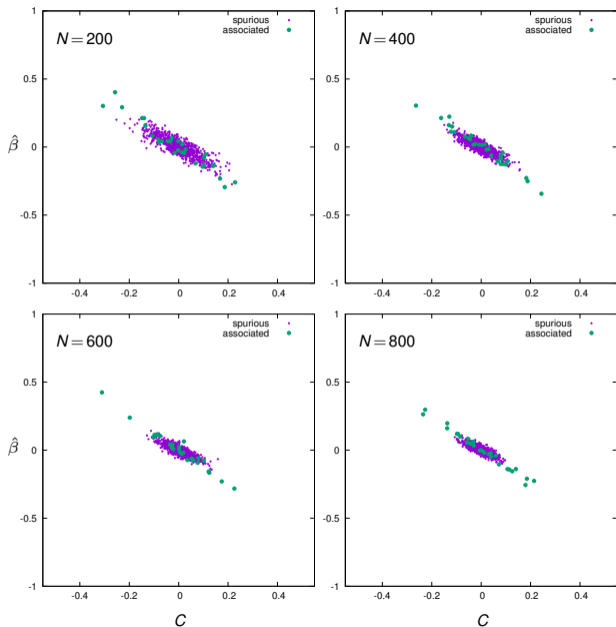
- ▶ $N \uparrow$: prob of false positive associations \downarrow
- ▶ $p \uparrow$: prob of false positive associations \uparrow

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 $p = 800$



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

Other extensions

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 $\theta_{\text{ML}} \neq \theta^*$ ✗

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 $\theta_{\text{ML}} \neq \theta^*$ ✗

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 $\theta_{\text{ML}} \neq \theta^*$ ✗

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

Other extensions

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*

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

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*

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

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

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

► *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:

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

interpretation:

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

Step 2 – remove minimisation over θ

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

► *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:

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

interpretation:

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

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

Other extensions

Summary

The replica method

aim: make hard analytical calculations easy ...

here: compute the average over \mathcal{D}

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

▶ 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

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

The replica method

aim: make hard analytical calculations easy ...

here: compute the average over \mathcal{D}

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

► 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

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

The replica method

aim: make hard analytical calculations easy ...

here: compute the average over \mathcal{D}

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

► 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

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

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

Other extensions

Summary

Application to Cox regression

- ▶ 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 : Lambert function

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

Application to Cox regression

- ▶ 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 : Lambert function

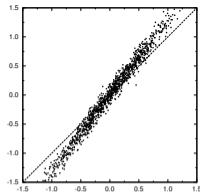
$$p(t|\xi, \lambda) = \lambda(t) e^{\xi - \exp(\xi)\Lambda(t)}$$

- ▶ interpretation:

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

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

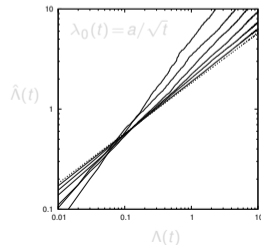
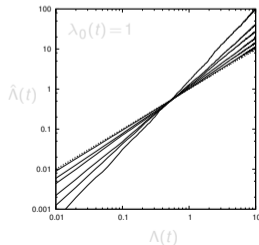
*all we need for
overfitting correction!*



- ▶ 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:

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

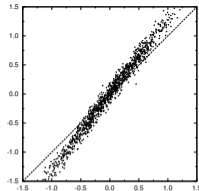
$\Lambda^*(t)$ drops out of equations!

- ▶ interpretation:

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

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

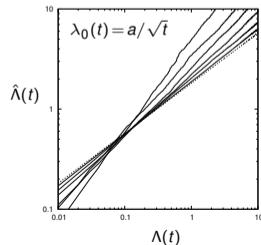
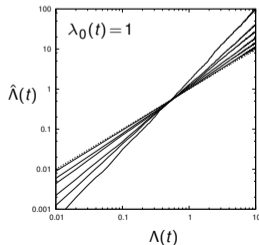
*all we need for
overfitting correction!*



- ▶ 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!

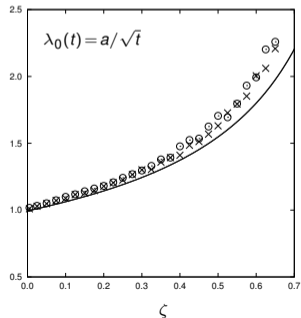
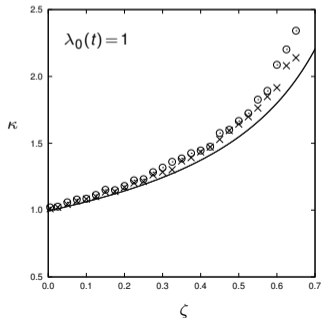
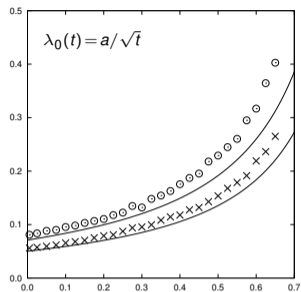
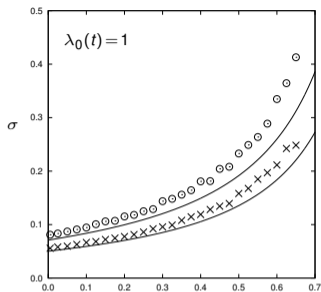
width σ and slope κ
of data clouds

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

simulations:

o: $N = 200$

x: $N = 400$



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

Other extensions

Summary

Application to regularized Cox regression

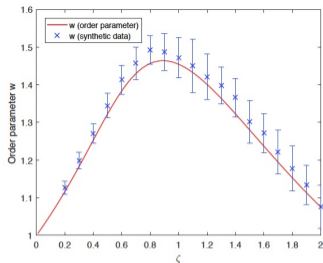
$N, p \rightarrow \infty, \zeta = p/N,$

L2-prior: $p(\beta) \propto \exp(-\eta p \beta^2)$

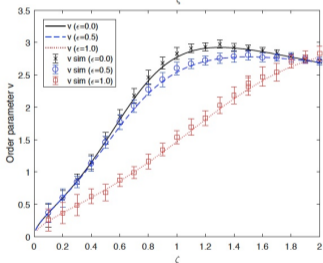
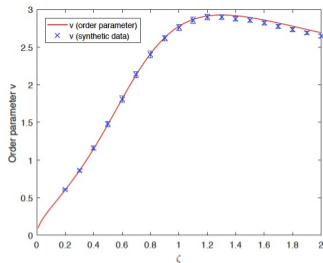
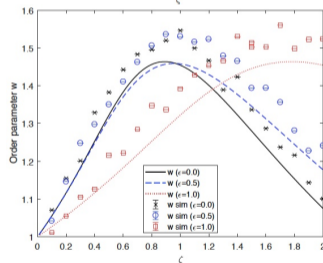
main changes:

- ▶ 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 $\kappa = 1$

*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 $\sigma = v$

final overfitting correction protocol

1. estimate covariate correlation matrix \mathbf{A} ,
compute its eigenvalue spectrum $\varrho(\mathbf{a})$
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

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

Other extensions

Summary

Other extensions

- ▶ survival analysis with censoring
 - end-of-trial censoring: minimal change to the theory
 - censoring by competing risks: more complex but doable
- ▶ generalized linear models:

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

e.g.

Logistic regression

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

Ordinal class regression

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

Latent class survival analysis

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

Other extensions

- ▶ survival analysis with censoring
 - end-of-trial censoring: minimal change to the theory
 - censoring by competing risks: more complex but doable
- ▶ generalized linear models:

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

e.g.

Logistic regression

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

Ordinal class regression

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

Latent class survival analysis

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

Summary

- ▶ *Overfitting in Cox regression*

- (i) bias in regression parameters: $\hat{\beta} \approx \kappa\beta^*$, $\kappa > 1$
- (ii) bias in base hazard rates
- (iii) extra noise, not captured by p -values

- ▶ *Analytical approach based on the replica method*

nontrivial closed equations for $\{u, v, w, \lambda(t)\}$

variational approximation for $\lambda(t)$

predictions for slope κ and noise σ

phase transition at $p/N = 1$

- ▶ *L2-regularized Cox regression*

theory involves spectrum of covariate correlation matrix

p -dependent hyperparameter

phase transition removed

reliable basis for overfitting corrections,
easily extended to arbitrary generalized linear models

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

Thank you!

Summary

- ▶ *Overfitting in Cox regression*

- (i) bias in regression parameters: $\hat{\beta} \approx \kappa\beta^*$, $\kappa > 1$
- (ii) bias in base hazard rates
- (iii) extra noise, not captured by p -values

- ▶ *Analytical approach based on the replica method*

nontrivial closed equations for $\{u, v, w, \lambda(t)\}$

variational approximation for $\lambda(t)$

predictions for slope κ and noise σ

phase transition at $p/N = 1$

- ▶ *L2-regularized Cox regression*

theory involves spectrum of covariate correlation matrix

p -dependent hyperparameter

phase transition removed

reliable basis for overfitting corrections,
easily extended to arbitrary generalized linear models

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

Thank you!

Summary

- ▶ *Overfitting in Cox regression*

- (i) bias in regression parameters: $\hat{\beta} \approx \kappa\beta^*$, $\kappa > 1$
- (ii) bias in base hazard rates
- (iii) extra noise, not captured by p -values

- ▶ *Analytical approach based on the replica method*

nontrivial closed equations for $\{u, v, w, \lambda(t)\}$

variational approximation for $\lambda(t)$

predictions for slope κ and noise σ

phase transition at $p/N = 1$

- ▶ *L2-regularized Cox regression*

theory involves spectrum of covariate correlation matrix

p -dependent hyperparameter

phase transition removed

reliable basis for overfitting corrections,
easily extended to arbitrary generalized linear models

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

Thank you!

Summary

- ▶ *Overfitting in Cox regression*

- (i) bias in regression parameters: $\hat{\beta} \approx \kappa\beta^*$, $\kappa > 1$
- (ii) bias in base hazard rates
- (iii) extra noise, not captured by p -values

- ▶ *Analytical approach based on the replica method*

nontrivial closed equations for $\{u, v, w, \lambda(t)\}$

variational approximation for $\lambda(t)$

predictions for slope κ and noise σ

phase transition at $p/N = 1$

- ▶ *L2-regularized Cox regression*

theory involves spectrum of covariate correlation matrix

p -dependent hyperparameter

phase transition removed

reliable basis for overfitting corrections,
easily extended to arbitrary generalized linear models

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

Thank you!