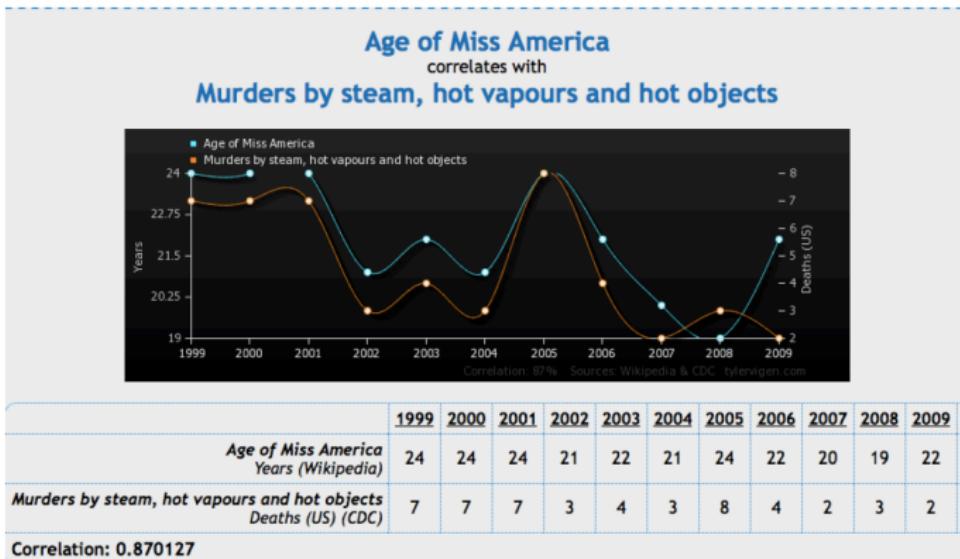


Replica analysis of overfitting in Cox regression

Statistical Physics Approaches to Systems Biology, Havana, Feb 2019

ACC Coolen

King's College London and Saddle Point Science



Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Proportional hazards regression (David Cox, 1972)

one risk,
no censoring:

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

parameters: $\beta, \lambda(t)$

- ▶ Maximum likelihood:

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

- ▶ Maximise over $\lambda(t)$ first

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

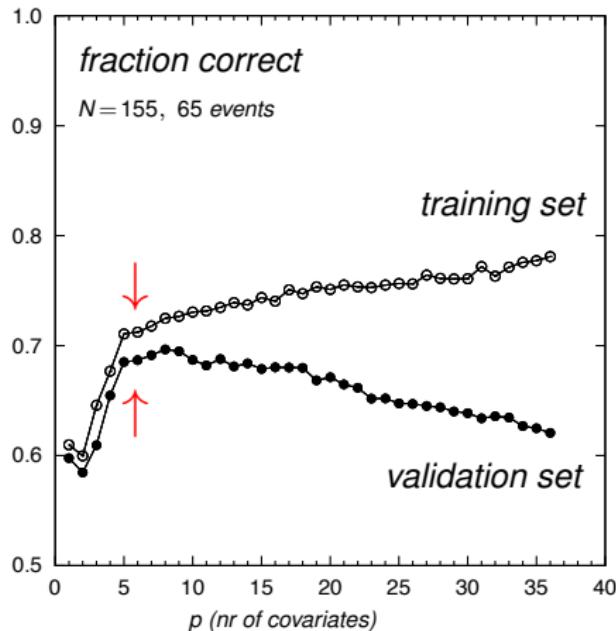
$$\hat{\beta} = \operatorname{argmax}_{\beta} \left\{ \sum_i \beta \cdot \mathbf{z}_i - \sum_i \log \left[\frac{\sum_j e^{\beta \cdot \mathbf{z}_j} \theta(t_j-t_i)}{\sum_j \theta(t_j-t_i)} \right] \right\}$$

overfitting in Cox regression

rule of thumb:

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

- ▶ too optimistic?
- ▶ dependent on β ?
- ▶ covariate correlations?



what happens in overfitting regime?

can we predict the optimal point?

analytical theory of overfitting in Cox regression?

theory simplification: choose $\langle z_\mu \rangle = 0$ and $\langle z_\mu z_\nu \rangle = \delta_{\mu\nu}$,
since average and correlations can be transformed away ...

- ▶ define

$$\langle f(\mathbf{z}) \rangle = \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i), \quad A_{\mu\nu} = \langle z_\mu z_\nu \rangle - \langle z_\mu \rangle \langle z_\nu \rangle, \quad \mathbf{z}_i = \langle \mathbf{z} \rangle + \mathbf{A}^{\frac{1}{2}} \tilde{\mathbf{z}}_i$$

- ▶ Cox regression:

$$\hat{\lambda}(t) = e^{-\hat{\beta} \cdot \langle \mathbf{z} \rangle} \frac{\sum_i \delta(t - t_i)}{\sum_i \theta(t_i - t) e^{\hat{\beta} \cdot \tilde{\mathbf{z}}_i}}$$

$$\hat{\beta} = \operatorname{argmax}_{\beta} \sum_i \left\{ (\mathbf{A}^{\frac{1}{2}} \beta) \cdot \tilde{\mathbf{z}}_i - \log \left[\sum_j \theta(t_j - t_i) e^{(\mathbf{A}^{\frac{1}{2}} \beta) \cdot \tilde{\mathbf{z}}_j} \right] \right\}$$

- ▶ put: $\hat{\beta} = \mathbf{A}^{-\frac{1}{2}} \tilde{\beta}$, $\hat{\lambda}(t) = \tilde{\lambda}(t) e^{-\tilde{\beta} \cdot \mathbf{A}^{-\frac{1}{2}} \bar{\mathbf{z}}}$

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

(standard eqns for $\{\tilde{\mathbf{z}}_i\}$)

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

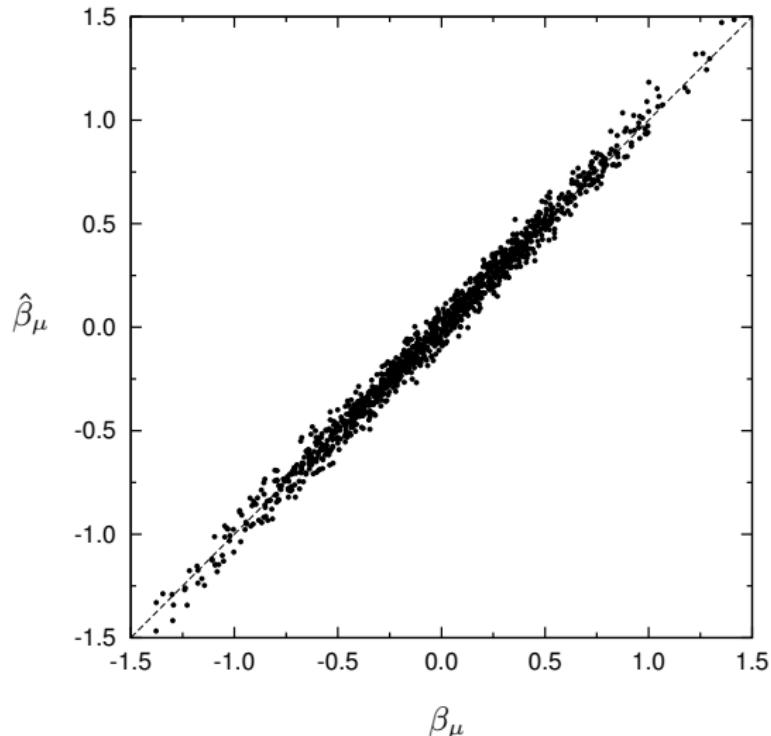
Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

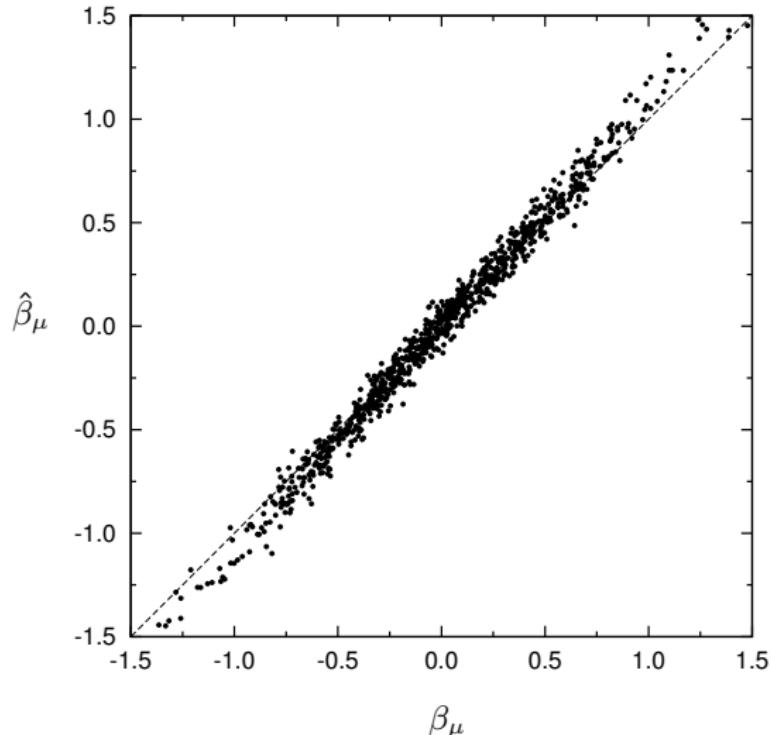
$N = 500$,
predicted versus true regression coefficients
(synthetic data, no censoring)

$$p/N = 0.002$$



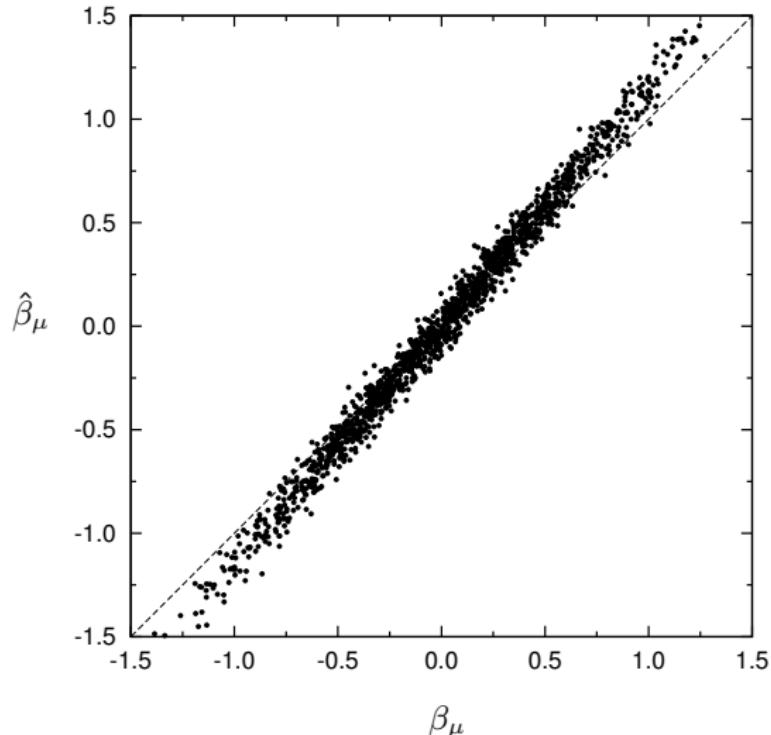
$N = 500$,
predicted versus true regression coefficients
(synthetic data, no censoring)

$$p/N = 0.10$$



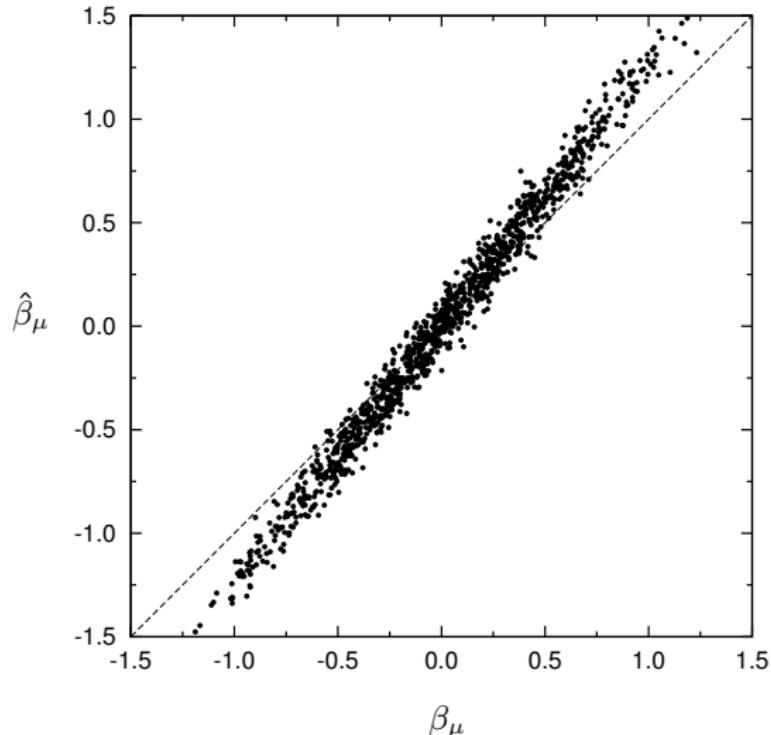
$N = 500$,
predicted versus true regression coefficients
(synthetic data, no censoring)

$$p/N = 0.20$$



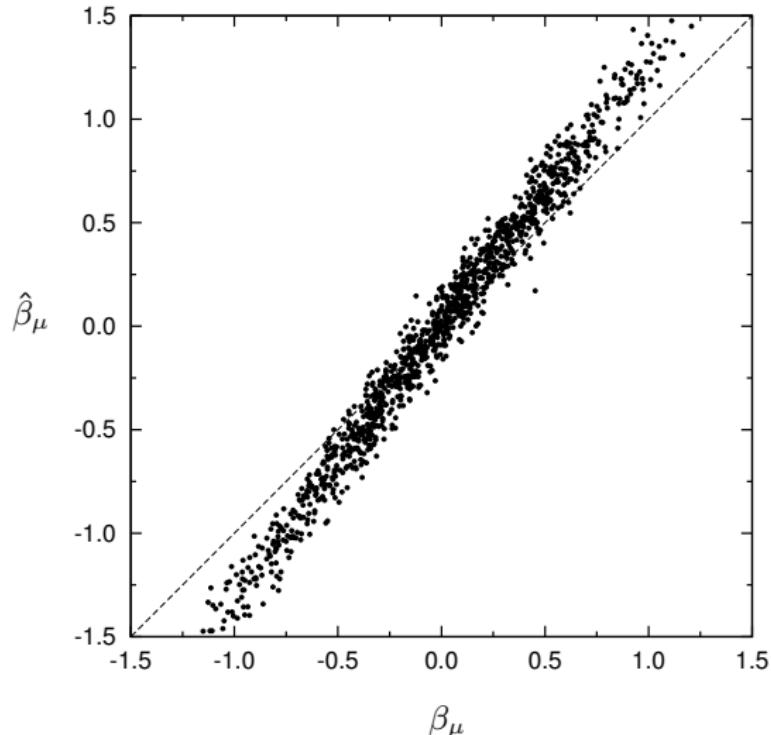
$N = 500$,
predicted versus true regression coefficients
(synthetic data, no censoring)

$$p/N = 0.30$$



$N = 500$,
predicted versus true regression coefficients
(synthetic data, no censoring)

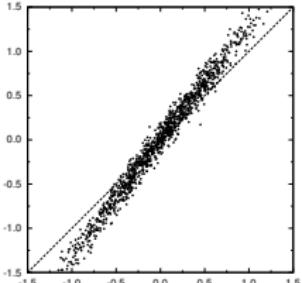
$$p/N = 0.40$$



Bad news

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

*we always inflate associations
(whether positive or negative)*

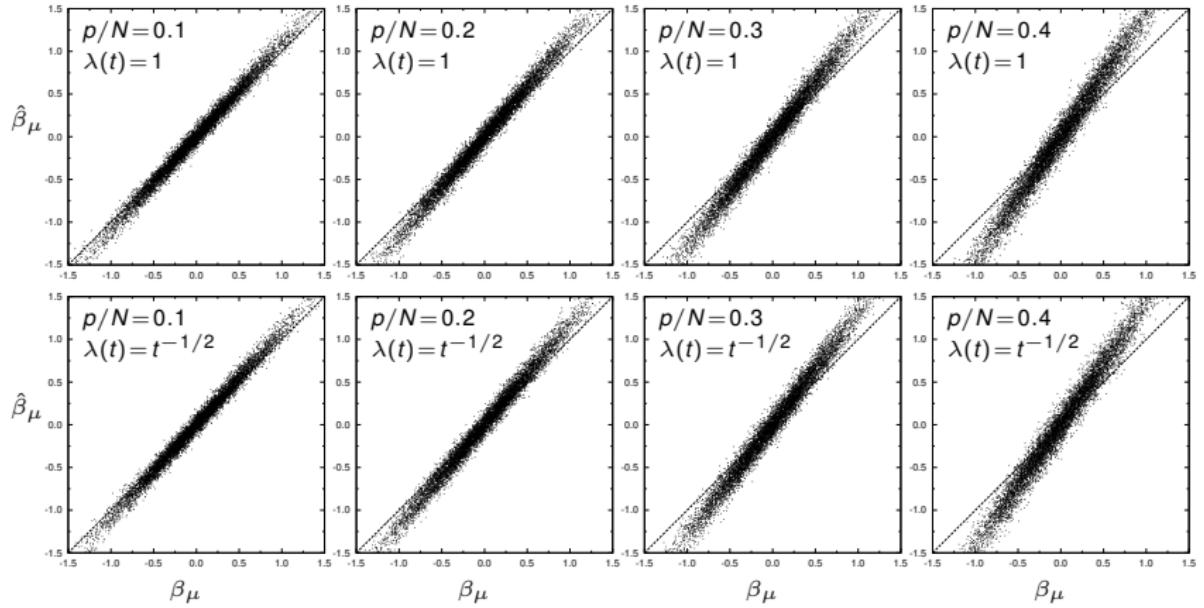


Good news

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

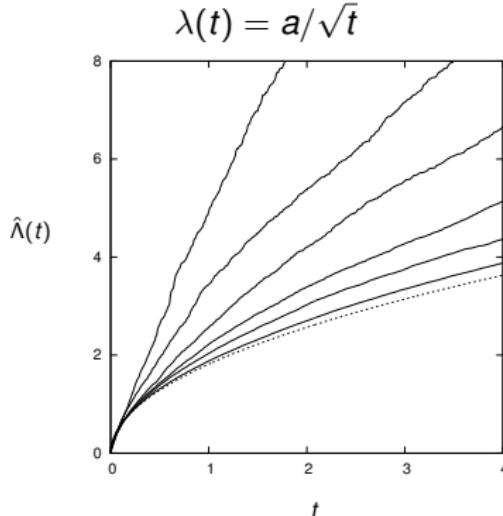
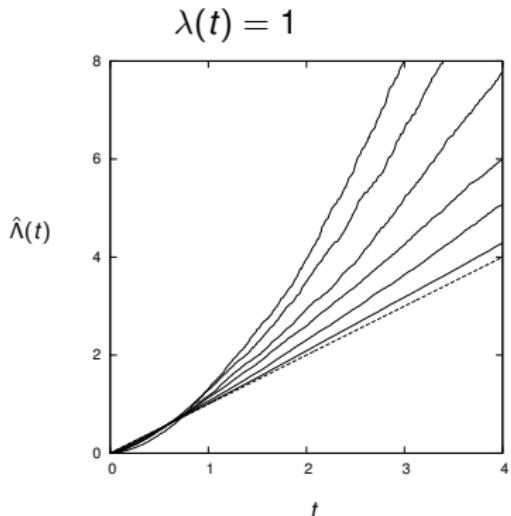
roadmap for research ...

- ▶ Predict impact of overfitting, in terms of
 - ratio p/N
 - covariate correlations
 - true parameters $\{\beta, \lambda(t)\}$
- ▶ Overfitting *correction* of Cox parameters
 - reliable regression at ratios $p/N \sim 0.5$ or more?



Association ‘inflation’ independent of true base hazard rate ...

$N = 400$,
 Gaussian association pars,
 $\langle \beta_\mu^2 \rangle = 0.25$



Base hazard rates underestimated for short times,
and over-estimated for large times ...

$$p/N = 0.05, 0.15, 0.25, 0.35, 0.45, 0.55
(\text{lower to upper curves})$$

Gaussian association pars, $\langle \beta_\mu^2 \rangle = 0.25$,
 $N = 400$, average event time $\langle t \rangle = 1$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Intuition for the problem ...

- ▶ *ML regression*

assumed model: $p_{\theta} = \{p(t|\mathbf{z}, \theta)\}$

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

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

ML pushes $p(t|\mathbf{z}, \theta)$ 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 $\theta_{\text{ML}} \neq \theta^*$ ✗

- ▶ *Barrier to overfitting theory*

want: relation between θ_{ML} and θ^* , for $p = \mathcal{O}(N)$

need: formula for θ_{ML} ...

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Step 1 – what to compute?

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

note:

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

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

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

$E(\theta^*, \mathcal{D}) = 0$: 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 – 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:

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

□

interpretation:

stochastic minimisation, with noise $\sim 1/\gamma$
(ground state = zero temp limit of free energy)

Step 3 – 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*

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

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

- ▶ *Application*

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

$$E_\gamma(\theta^*) = - \lim_{n \rightarrow 0} \frac{1}{Nn} \frac{\partial}{\partial \gamma} \log \int d\theta^1 \dots d\theta^n \left[\int d\mathbf{z} dt p(\mathbf{z}) p(t | \mathbf{z}, \theta^*) \prod_{\alpha=1}^n \left[\frac{p(t | \mathbf{z}, \theta^\alpha)}{p(t | \mathbf{z}, \theta^*)} \right]^\gamma \right]^N$$

□

still completely general ...

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

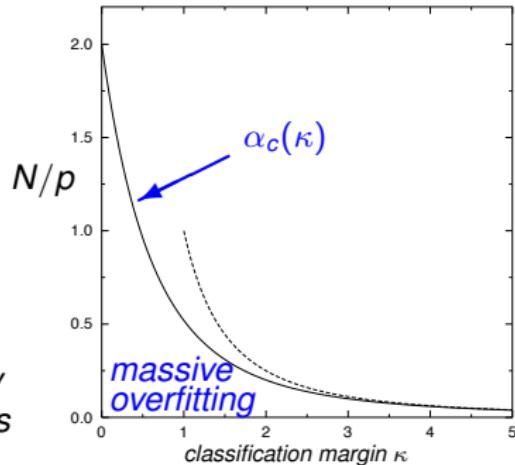
heterogeneous stochastic systems in physics,
biology, computer science, economics, ...

- ▶ 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

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Translation to Cox's model

- Here $\theta = \{\lambda, \beta\}$, and

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

$$\begin{aligned} E_\gamma(\beta^*, \lambda^*) &= -\frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{Nn} \log \int \{d\lambda_1 \dots d\lambda_n\} \int d\beta^1 \dots d\beta^n \\ &\times \left\{ \int d\mathbf{z} dt p(\mathbf{z}) p(t|\mathbf{z}, \lambda^*, \beta^*) \prod_{\alpha=1}^n \left[\frac{p(t|\mathbf{z}, \lambda_\alpha, \beta^\alpha)}{p(t|\mathbf{z}, \lambda^*, \beta^*)} \right]^\gamma \right\}^N \end{aligned}$$

- short-hands:

$$p(t|\xi, \lambda) = \lambda(t) e^{\xi - \exp(\xi) \int_0^t ds \lambda(s)}, \quad (\beta^*, \lambda^*) = (\beta^0, \lambda_0)$$

$$\mathbf{y} = (y_0, \dots, y_p), \quad p(\mathbf{y}|\beta^0, \dots, \beta^n) = \int d\mathbf{z} p(\mathbf{z}) \prod_{\alpha=0}^n \delta \left[y_\alpha - \frac{\beta^\alpha \cdot \mathbf{z}}{\sqrt{p}} \right]$$

$$\begin{aligned} E_\gamma(\beta^0, \lambda_0) &= -\frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{Nn} \log \int \{d\lambda_1 \dots d\lambda_n\} \int d\beta^1 \dots d\beta^n \\ &\times \left\{ \int d\mathbf{y} p(\mathbf{y}|\beta^0, \dots, \beta^n) \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left[\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right]^\gamma \right\}^N \end{aligned}$$

- ▶ all $\{y_\alpha\}$: sums of zero-average Gaussian random vars, so

$$p(\mathbf{y}|\boldsymbol{\beta}^0, \dots, \boldsymbol{\beta}^n) = \frac{e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{C}^{-1}[\{\boldsymbol{\beta}\}] \mathbf{y}}}{\sqrt{(2\pi)^{n+1} \text{Det} \mathbf{C}[\{\boldsymbol{\beta}\}]}}$$

$$C_{\alpha\rho}[\{\boldsymbol{\beta}\}] = \frac{1}{p} \int d\mathbf{z} p(\mathbf{z}) (\boldsymbol{\beta}^\alpha \cdot \mathbf{z}) (\boldsymbol{\beta}^\rho \cdot \mathbf{z}) = \frac{1}{p} \boldsymbol{\beta}^\alpha \cdot \boldsymbol{\beta}^\rho$$

- ▶ introduce for each (α, ρ) :

$$1 = \int dC_{\alpha\rho} \delta[C_{\alpha\rho} - C_{\alpha\rho}[\{\boldsymbol{\beta}\}]] = \int \frac{dC_{\alpha\rho} d\hat{C}_{\alpha\rho}}{2\pi/p} e^{ip\hat{C}_{\alpha\rho} [C_{\alpha\rho} - C_{\alpha\rho}[\{\boldsymbol{\beta}\}]]}$$

after simple manipulations,

$$\begin{aligned} E_\gamma(\boldsymbol{\beta}^0, \lambda_0) &= -\frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{Nn} \log \int \{d\lambda_1 \dots d\lambda_n\} \int \frac{d\mathbf{C} d\hat{\mathbf{C}} e^{ip \sum_{\alpha\rho=0}^n \hat{C}_{\alpha\rho} C_{\alpha\rho}}}{(2\pi/p)^{(n+1)^2}} \\ &\times \left\{ \int \frac{d\mathbf{y} e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}}}{\sqrt{(2\pi)^{n+1} \text{Det} \mathbf{C}}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left[\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right]^\gamma \right\}^N \\ &\times \int d\boldsymbol{\beta}^1 \dots d\boldsymbol{\beta}^n e^{-i \sum_{\alpha\rho=0}^n \hat{C}_{\alpha\rho} \boldsymbol{\beta}^\alpha \cdot \boldsymbol{\beta}^\rho} \end{aligned}$$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Conversion to saddle-point problem

- ▶ define: $P(\beta_0) = \frac{1}{p} \sum_{\mu=1}^p \delta[\beta_0 - \beta_\mu^0]$, $S^2 = \int d\beta_0 P(\beta_0) \beta_0^2$

$$E_\gamma(P, \lambda_0) = -\frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{N^n} \log \int \{d\lambda_1 \dots d\lambda_n\} \int \frac{d\mathbf{C} d\hat{\mathbf{C}} e^{ip(\sum_{\alpha\rho=0}^n \hat{c}_{\alpha\rho} c_{\alpha\rho} - \hat{c}_{00} S^2)}}{(2\pi/p)^{(n+1)^2}}$$

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

$$\times e^{p \int d\beta_0 P(\beta_0) \log \int d\beta_1 \dots d\beta_n e^{-2i\beta_0 \sum_{\rho=1}^n \hat{c}_{0\rho} \beta_\rho - i \sum_{\alpha\rho=1}^n \hat{c}_{\alpha\rho} \beta_\alpha \beta_\rho}}$$

- ▶ exchange order of $n \rightarrow 0$ and $N \rightarrow \infty$,
with fixed $p/N = \zeta$: steepest descent integral ...

$$\lim_{N \rightarrow \infty} E_\gamma(P, \lambda_0) = \frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{n} \text{extr}_{\mathbf{C}, \hat{\mathbf{C}}, \lambda_1, \dots, \lambda_n} \Psi[\mathbf{C}, \hat{\mathbf{C}}; \lambda_1, \dots, \lambda_n]$$

$$\Psi[\dots] = -i\zeta \left[\sum_{\alpha\rho=0}^n \hat{c}_{\alpha\rho} c_{\alpha\rho} - \hat{c}_{00} S^2 \right] + \frac{1}{2}(n+1) \log(2\pi) + \frac{1}{2} \log \text{Det} \mathbf{C}$$

$$- \zeta \int d\beta_0 P(\beta_0) \log \int d\beta_1 \dots d\beta_n e^{-2i\beta_0 \sum_{\rho=1}^n \hat{c}_{0\rho} \beta_\rho - i \sum_{\alpha\rho=1}^n \hat{c}_{\alpha\rho} \beta_\alpha \beta_\rho}$$

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

- ▶ simple steps ...

$$C_{00} = S^2, \quad \alpha, \rho = 1 \dots n : \quad \hat{C}_{\alpha\rho} = -\frac{1}{2}iD_{\alpha\rho}, \quad \hat{C}_{0\rho} = -\frac{1}{2}id_\rho$$

with $D_{\alpha\rho}, d_\rho \in \mathbb{R}$, and $\mathbf{D} = \{D_{\alpha\rho}\}$ positive definite,

$$\begin{aligned} \Psi[\dots] &= -\frac{1}{2}\zeta \sum_{\alpha\rho=1}^n D_{\alpha\rho} C_{\alpha\rho} - \zeta \sum_{\rho=1}^n d_\rho C_{0\rho} - \frac{1}{2}\zeta S^2 \sum_{\alpha\rho=1}^n d_\alpha (\mathbf{D}^{-1})_{\alpha\rho} d_\rho \\ &\quad + \frac{1}{2}(n+1) \log(2\pi) + \frac{1}{2} \log \text{Det} \mathbf{C} - \zeta \log \int d\beta_1 \dots d\beta_n e^{-\frac{1}{2} \sum_{\alpha\rho=1}^n D_{\alpha\rho} \beta_\alpha \beta_\rho} \\ &\quad - \log \int d\mathbf{y} e^{-\frac{1}{2} \mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left[\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right]^\gamma \end{aligned}$$

- ▶ variation with respect to $\{d_\alpha\}$: $d_\alpha = -\sum_\rho D_{\alpha\rho} C_{0\rho} / S^2$

$$\begin{aligned} \Psi[\dots] &= -\frac{1}{2}\zeta \sum_{\alpha\rho=1}^n D_{\alpha\rho} \left[C_{\alpha\rho} - \frac{C_{0\alpha} C_{0\rho}}{S^2} \right] + \frac{1}{2}(n+1) \log(2\pi) + \frac{1}{2} \log \text{Det} \mathbf{C} \\ &\quad - \log \int d\mathbf{y} e^{-\frac{1}{2} \mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left[\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right]^\gamma \\ &\quad - \zeta \log \int d\beta_1 \dots d\beta_n e^{-\frac{1}{2} \sum_{\alpha\rho=1}^n D_{\alpha\rho} \beta_\alpha \beta_\rho} \end{aligned}$$

- ▶ eqns depend on $P(\beta_0)$ only via $S^2 = \int d\beta_0 P(\beta_0)\beta_0^2$
- ▶ variation with respect to \mathbf{D} : $(\mathbf{D}^{-1})_{\alpha\rho} = C_{\alpha\rho} - C_{0\alpha}C_{0\rho}/S^2$,

$$E_\gamma(S, \lambda_0) = \frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{n} \text{extr}_{\mathbf{C}; \lambda_1, \dots, \lambda_n} \Psi[\mathbf{C}; \lambda_1, \dots, \lambda_n]$$

$$\begin{aligned} \Psi[\mathbf{C}; \lambda_1, \dots, \lambda_n] &= \frac{1}{2} \log \text{Det} \mathbf{C} - \frac{1}{2} \zeta \log \text{Det} \mathbf{C}' \\ &\quad - \log \int \frac{d\mathbf{y}}{\sqrt{2\pi}} e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left[\frac{p(t|y_\alpha, \lambda_\alpha)}{p(t|y_0, \lambda_0)} \right]^\gamma \end{aligned}$$

notes:

modulo terms that vanish due to $n \rightarrow 0$ or $\partial/\partial\gamma$,

extremisation over \mathbf{C} : subject to $C_{00} = S^2$,

\mathbf{C}' : $n \times n$ matrix, $C'_{\alpha\rho} = C_{\alpha\rho} - C_{0\alpha}C_{0\rho}/S^2$ ($\alpha, \rho = 1 \dots n$)

*to proceed:
ansatz for form of $\{C_{\alpha\beta}, \lambda_\alpha\}$*

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Replica symmetric saddle points

solution space of the regression algorithm *ergodic*,
set of equivalent minima in regression parameter space *connected*

$$\forall \alpha, \rho = 1 \dots n : \quad \lambda_\alpha(t) = \lambda(t), \quad C_{0\alpha} = c_0,$$
$$C_{\alpha\rho} = C\delta_{\alpha\rho} + c(1 - \delta_{\alpha\rho})$$

$$E_\gamma(S, \lambda_0) = \frac{\partial}{\partial \gamma} \text{extr}_{C, c, c_0; \lambda} \Psi_{\text{RS}}[C, c, c_0; \lambda]$$

$$\begin{aligned} \Psi_{\text{RS}}[C, c, c_0; \lambda] &= \lim_{n \rightarrow 0} \frac{1}{n} \left\{ \frac{1}{2} \log \text{Det} \mathbf{C} - \frac{1}{2} \zeta \log \text{Det} \mathbf{C}' \right. \\ &\quad \left. - \log \int \frac{d\mathbf{y}}{\sqrt{2\pi}} e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{C}^{-1} \mathbf{y}} \int dt p(t|y_0, \lambda_0) \prod_{\alpha=1}^n \left[\frac{p(t|y_\alpha, \lambda)}{p(t|y_0, \lambda_0)} \right]^\gamma \right\} \end{aligned}$$

$$\mathbf{C} = \begin{pmatrix} S^2 & c_0 & \cdots & \cdots & c_0 \\ c_0 & C & c & \cdots & c \\ \vdots & c & C & \cdots & c \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_0 & c & \cdots & c & C \end{pmatrix}$$

- eigenvectors and eigenvalues of \mathbf{C} :

$$(u, v, \dots, v) : \quad \mu_{\pm} = \frac{1}{2} \left\{ C + (n-1)c + S^2 \pm \sqrt{[C + (n-1)c - S^2]^2 + 4nc_0^2} \right\}$$

$$(0, v_1, \dots, v_n) : \quad \sum_{\alpha=1}^n v_{\alpha} = 0, \quad \mu = C - c \quad (\text{multiplicity } n-1)$$

so

$$\begin{aligned} \log \text{Det} \mathbf{C} &= \log [(C - c)^{n-1} \mu_+ \mu_-] \\ &= \log S^2 + n \log(C - c) + n \frac{c - c_0^2 / S^2}{C - c} + \mathcal{O}(n^2) \end{aligned}$$

- form of \mathbf{C}^{-1} :

$$\mathbf{C}^{-1} = \begin{pmatrix} d_{00} & d_0 & \cdots & \cdots & d_0 \\ d_0 & D & d & \cdots & d \\ \vdots & d & D & \cdots & d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ d_0 & d & \cdots & d & D \end{pmatrix}$$

$$d_{00} = \frac{C + (n-1)c}{S^2(C + (n-1)c) - nc_0^2}, \quad d_0 = -\frac{c_0}{S^2(C + (n-1)c) - nc_0^2}$$

$$d = \frac{1}{C - c} \frac{c_0^2 - cS^2}{S^2(C + (n-1)c) - nc_0^2}, \quad D = d + \frac{1}{C - c}$$

- eigenvectors and eigenvalues of \mathbf{C}' :

$$(1, \dots, 1) : \quad \text{eigenvalue } C - c - nc_0^2/S^2 + nc$$

$$(1, \dots, 1)^\perp : \quad \text{eigenvalue } C - c \quad (n-1 \text{ fold degenerate})$$

so

$$\begin{aligned} \log \text{Det} \mathbf{C}' &= (n-1) \log(C-c) + \log \left(C - c + n(c - c_0^2/S^2) \right) \\ &= n \left[\log(C-c) + (c - c_0^2/S^2)(C-c) \right] + \mathcal{O}(n^2) \end{aligned}$$

- insert into Ψ_{RS} ,
expand for small n :

$$\begin{aligned} \Psi_{\text{RS}}[C, c, c_0; \lambda] &= \frac{1}{2}(1-\zeta) \left[\log(C-c) + \frac{c - c_0^2/S^2}{C-c} \right] + \frac{1}{2} \frac{c_0^2/S^2}{C-c} - \frac{1}{2} \log(C-c) \\ &\quad - \frac{1}{2} \log(2\pi) - \int Dz Dy_0 \int dt p(t|Sy_0, \lambda_0) \\ &\quad \times \log \int Dy e^{y[y_0 c_0 / S \sqrt{C-c} + z \sqrt{(c - c_0^2/S^2)/(C-c)}]} \left(\frac{p(t|y \sqrt{C-c}, \lambda)}{p(t|Sy_0, \lambda_0)} \right)^\gamma \end{aligned}$$

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

- ▶ transform RS order pars: $u = \sqrt{C - c}$, $v = \sqrt{c - c_0^2/S^2}$, $w = c_0/S$,
 transform integration var: $y \rightarrow y + (wy_0 + vz)/u$,

$$\begin{aligned}
 E_\gamma(S, \lambda_0) &= \int Dy_0 \int dt p(t|Sy_0, \lambda_0) \log p(t|Sy_0, \lambda_0) \\
 &\quad - \frac{\partial}{\partial \gamma} \text{extr}_{u,v,w;\lambda} \left\{ \zeta \left(\frac{v^2}{2u^2} + \log u \right) \right. \\
 &\quad \left. + \int Dz Dy_0 \int dt p(t|Sy_0, \lambda_0) \log \int Dy p^\gamma(t|uy + wy_0 + vz, \lambda) \right\} \\
 &= \int Dy_0 \int dt p(t|Sy_0, \lambda_0) \left\{ \log p(t|Sy_0, \lambda_0) \right. \\
 &\quad \left. - \int Dz \left[\frac{\int Dy p^\gamma(t|uy + wy_0 + vz, \lambda) \log p(t|uy + wy_0 + vz, \lambda)}{\int Dy p^\gamma(t|uy + wy_0 + vz, \lambda)} \right] \right\}
 \end{aligned}$$

$\{u, v, w; \lambda\}$: evaluated at saddle point of

$$\begin{aligned}
 \Psi_{\text{RS}}(u, v, w; \lambda) &= \zeta \left(\frac{v^2}{2u^2} + \log u \right) \\
 &\quad + \int Dz Dy_0 \int dt p(t|Sy_0, \lambda_0) \log \int Dy p^\gamma(t|uy + wy_0 + vz, \lambda)
 \end{aligned}$$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Physical interpretation of order parameters

- ▶ usual replica manipulations:
(fraction form, or weak perturbations):

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

transformed order pars:

$$u^2 = \lim_{p \rightarrow \infty} \frac{1}{p} \langle \langle \beta^2 \rangle - \langle \beta \rangle^2 \rangle_{\mathcal{D}}, \quad w = \lim_{p \rightarrow \infty} \frac{1}{\sqrt{p}} \frac{\beta^* \cdot \langle \langle \beta \rangle \rangle_{\mathcal{D}}}{|\beta^*|}$$
$$v^2 = \lim_{p \rightarrow \infty} \frac{1}{p} \left[\langle \langle \beta \rangle^2 \rangle_{\mathcal{D}} - \left(\frac{\beta^* \cdot \langle \langle \beta \rangle \rangle_{\mathcal{D}}}{|\beta^*|} \right)^2 \right]$$

- ▶ unique ML point: $\lim_{\gamma \rightarrow \infty} u = 0$
(true for Cox regression, minimization of convex function)

perfect regression:

$\beta = \beta^*$ for all \mathcal{D}, β^* :

$$c_0 = c = C = S^2, \quad u = v = 0, \quad w = S$$

- ▶ if $\langle \beta \rangle \approx \kappa \beta^* + \xi$,
with $\langle \xi \rangle_{\mathcal{D}} = \mathbf{0}$, $\langle \xi^2 \rangle_{\mathcal{D}} / p = \sigma^2$:

$$\kappa = w/S, \quad \sigma = v \quad (\text{slope and width of parameter cloud})$$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

RS saddle point eqns

- eqns for (u, v, w) , using $\partial \log p(t|\xi) / \partial \xi = 1 - e^\xi \Lambda(t)$:

$$\frac{\zeta}{\gamma u^2} \left(\frac{v^2}{\gamma u^2} - \frac{1}{\gamma} \right) = \int Dz Dy_0 \int dt p(t|Sy_0, \lambda_0)$$

$$\times \frac{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda) \left[[1 - e^{uy+wy_0+vz} \Lambda(t)]^2 - \gamma^{-1} e^{uy+wy_0+vz} \Lambda(t) \right]}{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda)}$$

$$\frac{\zeta v}{\gamma u^2} = \int Dz Dy_0 z \int dt p(t|Sy_0, \lambda_0) \Lambda(t) \frac{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda) e^{uy+wy_0+vz}}{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda)}$$

$$0 = \int Dz Dy_0 y_0 \int dt p(t|Sy_0, \lambda_0) \Lambda(t) \frac{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda) e^{uy+wy_0+vz}}{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda)}$$

- functional order parameter eqn,
 $\delta \Psi_{\text{RS}}(u, v, w; \lambda) / \delta \lambda(s) = 0$:

$$\frac{p(s)}{\lambda(s)} = \int Dz Dy_0 \int_s^\infty dt p(t|Sy_0, \lambda_0) \frac{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda) e^{uy+wy_0+vz}}{\int Dy p^\gamma(t|uy+wy_0+vz, \lambda)}$$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

- correct scaling: $u = \tilde{u}/\sqrt{\gamma}$,
tricky manipulations to get $\gamma \rightarrow \infty$ limit:

$$\zeta v^2 = \int Dz Dy_0 \int dt p(t|Sy_0, \lambda_0) \left[\tilde{u}^2 - W(\tilde{u}^2 e^{\tilde{u}^2 + wy_0 + vz} \Lambda(t)) \right]^2$$

$$\zeta v = \int Dz Dy_0 z \int dt p(t|Sy_0, \lambda_0) W(\tilde{u}^2 e^{\tilde{u}^2 + wy_0 + vz} \Lambda(t))$$

$$0 = \int Dz Dy_0 y_0 \int dt p(t|Sy_0, \lambda_0) W(\tilde{u}^2 e^{\tilde{u}^2 + wy_0 + vz} \Lambda(t))$$

$$\frac{p(t)}{\lambda(t)} = \int Dz Dy_0 \int_t^\infty dt' p(t'|Sy_0, \lambda_0) \frac{W(\tilde{u}^2 e^{\tilde{u}^2 + wy_0 + vz} \Lambda(t'))}{\tilde{u}^2 \Lambda(t')}$$

$W(z)$: Lambert's W -function
(inverse of $f(x) = xe^x$)

- overfitting measure:

$$E(S, \lambda_0) = \int dt p(t) \log \left[\frac{\lambda_0(t)}{\lambda(t)} \right] - (1 + \tilde{u}^2) \left[1 - \frac{1}{\tilde{u}^2} \int Dz Dy_0 \int dt p(t|Sy_0, \lambda_0) W(\tilde{u}^2 e^{\tilde{u}^2 + wy_0 + vz} \Lambda(t)) \right]$$

- $\zeta \rightarrow 0$: correct soln, no overfitting ($\tilde{u} = v = 0, w = S$)
- $\zeta \rightarrow 1$: phase transition ($\tilde{u}, v, w \rightarrow \infty$)

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

preparations:

- ▶ write functional eqn in form involving $\Lambda(t)$ only:

$$\Lambda(t) = \int_0^t dt' p(t') \left\{ \int Dz Dy_0 \int_{t'}^\infty dt'' p(t'' | Sy_0, \lambda_0) \frac{W(\tilde{u}^2 e^{\tilde{u}^2 + wy_0 + vz} \Lambda(t''))}{\tilde{u}^2 \Lambda(t'')} \right\}^{-1}$$

- ▶ transform integration over t into
integration over survival function $s = \exp[-e^{Sy_0} \Lambda_0(t)]$,

$$\text{short-hand } L(t) = \tilde{u}^2 e^{\tilde{u}^2} \Lambda(t),$$

$$\zeta v^2 = \int Dy_0 Dz \int_0^1 ds \left[\tilde{u}^2 - W(e^{wy_0 + vz} L(t(s, y_0))) \right]^2$$

$$\zeta = \int Dy_0 Dz \int_0^1 ds \left\{ \frac{W(e^{wy_0 + vz} L(t(s, y_0)))}{1 + W(e^{wy_0 + vz} L(t(s, y_0)))} \right\}$$

$$\frac{\zeta w}{s} = - \int Dy_0 Dz \int_0^1 ds [1 + \log(s)] W(e^{wy_0 + vz} L(t(s, y_0)))$$

$$L(t) = \tilde{u}^2 \int_0^t dt' p(t') \left\{ \int Dy_0 Dz \int_0^1 ds' \frac{\theta[t(s', y_0) - t']}{L(t(s', y_0))} W(e^{wy_0 + vz} L(t(s', y_0))) \right\}^{-1}$$

- ▶ main mathematical challenge: functional eqn

write $L(t) = \Phi(\Lambda_0(t))$:

$$\frac{\tilde{u}^2 g(x)}{\frac{d \log \Phi(x)}{dx}} \frac{d}{dx} \log \left(\frac{d\Phi(x)/dx}{g(x)} \right) = \int D y_0 e^{S y_0 - x \exp(S y_0)} \int D z W(e^{w y_0 + v z} \Phi(x))$$

$$g(x) = \int D y e^{S y - x \exp(S y)}$$

- ▶ asymptotic form, $x \rightarrow \infty$, using

$$W(z) = \log z - \log(\log z) + \mathcal{O}(\log(\log z) / \log z) \quad (z \rightarrow \infty)$$

$$\log g(x) = -\frac{1}{2S^2} (\log x)^2 + \frac{1}{S^2} \log x \cdot \log(\log x) + \mathcal{O}(\log x) \quad (x \rightarrow \infty)$$

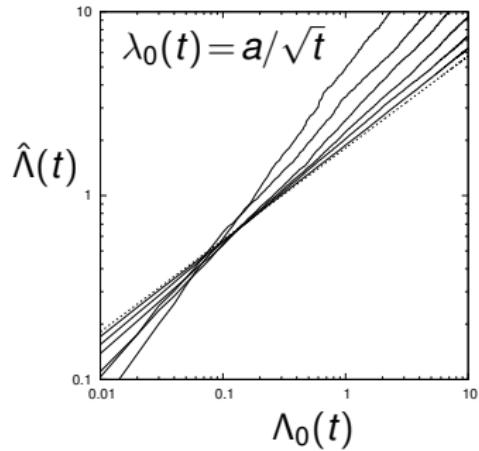
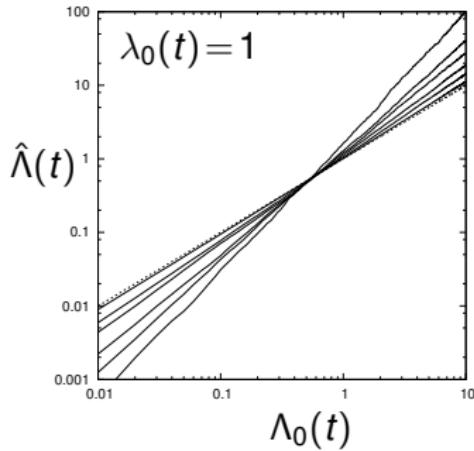
$$\Phi(x) = \rho \log x + (1-\rho) \log \log x + o(\log \log x)$$

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

hence

$$t \gg 1 : \quad \log \Lambda(t) = \rho \log \Lambda_0(t) + (1-\rho) \log(\log \Lambda_0(t)) + \dots$$

leading order : $\Lambda(t) \approx \Lambda_0^\rho(t)$



simulation results:

$$\zeta = 0.05, 0.15, 0.25, 0.35, 0.45, 0.55$$

confirms $\log \Lambda(t) \approx \rho \log \Lambda_0(t)$,
with $\rho = \rho(\zeta)$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Variational approximation

$$\Lambda(t) = k\Lambda_0^\rho(t),$$

insert into Ψ_{RS} , extremize over $(k, \rho, \tilde{u}, v, w)$

write $q = k\tilde{u}^2 e^{\tilde{u}^2}$,

further manipulations of integrals ...

$$\zeta v^2 = \int Dx \int_0^1 ds \left[\tilde{u}^2 - W\left(qe^{x\sigma(v,w,\rho)} \log^\rho(1/s)\right) \right]^2$$

$$\zeta = \int Dx \int_0^1 ds \frac{W\left(qe^{x\sigma(v,w,\rho)} \log^\rho(1/s)\right)}{1 + W\left(qe^{x\sigma(v,w,\rho)} \log^\rho(1/s)\right)}$$

$$\frac{\zeta w}{s} = - \int Dx \int_0^1 ds [1 + \log(s)] W\left(qe^{x\sigma(v,w,\rho)} \log^\rho(1/s)\right)$$

$$\tilde{u}^2 = \int Dx \int_0^1 ds W\left(qe^{x\sigma(v,w,\rho)} \log^\rho(1/s)\right)$$

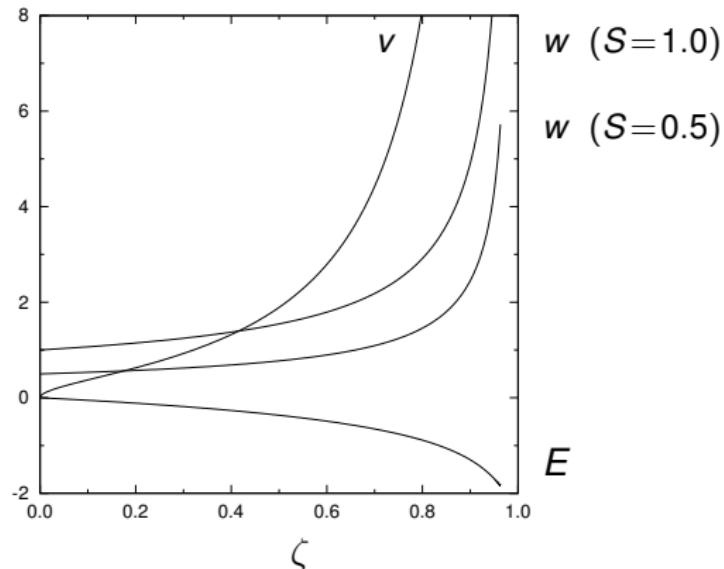
$$\frac{\tilde{u}^2}{\rho} = \int Dx \int_0^1 ds W\left(qe^{x\sigma(v,w,\rho)} \log^\rho(1/s)\right) \log \log(1/s) - S(w - \rho S) \zeta + \tilde{u}^2 C_E$$

C_E : Euler's constant

- ▶ similarly:

$$\begin{aligned} E(S, \lambda_0) &= \int dt p(t) \log \left[\lambda_0(t)/\lambda(t) \right] \\ &= -\log k - \log \rho + (\rho - 1) C_E \end{aligned}$$

- ▶ numerical soln
of coupled eqns
for $(k, \rho, \tilde{u}, v, w)$
completely indep
of true base hazard
rate $\lambda_0(t)$



Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

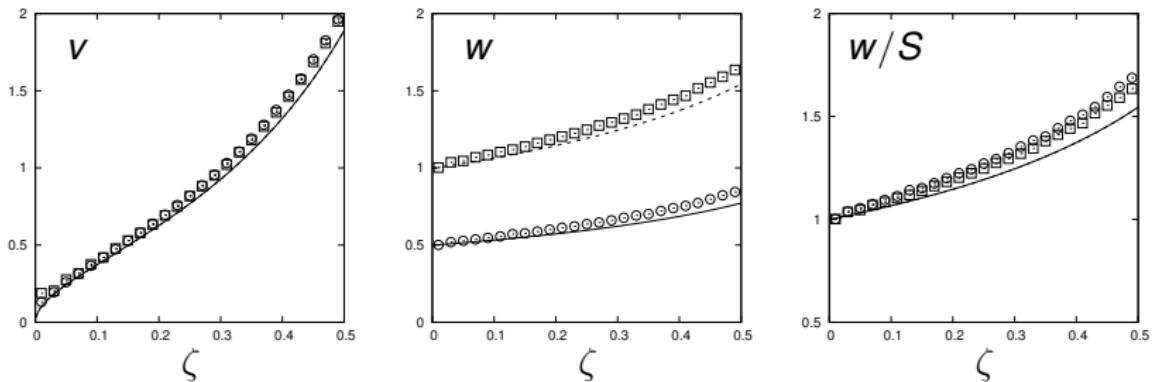
Summary

Beyond replica symmetry

Censoring, competing risks and priors

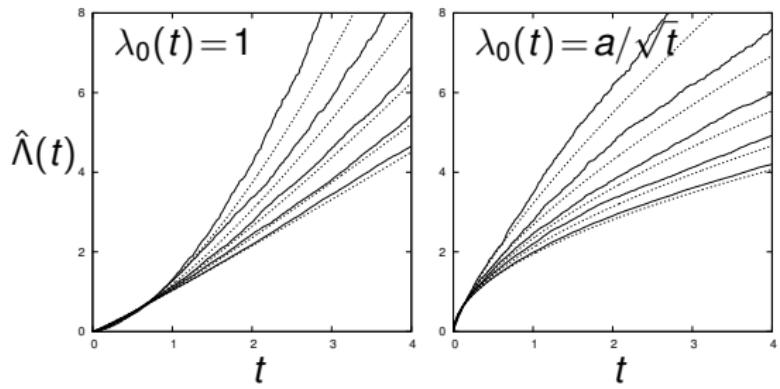
Overfitting in inference with max entropy models

Variational theory: tests of order parameters



simulations: $N=200$, $\lambda_0(t) = 1$, $S=0.5$ (circles) or $S=1$ (squares)

$N=400$, $S=0.5$,
 a such that $\int dt p(t)t = 1$,
 $\zeta = 0.1, 0.2, 0.3, 0.4, 0.5$
(lower to upper curves)



Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

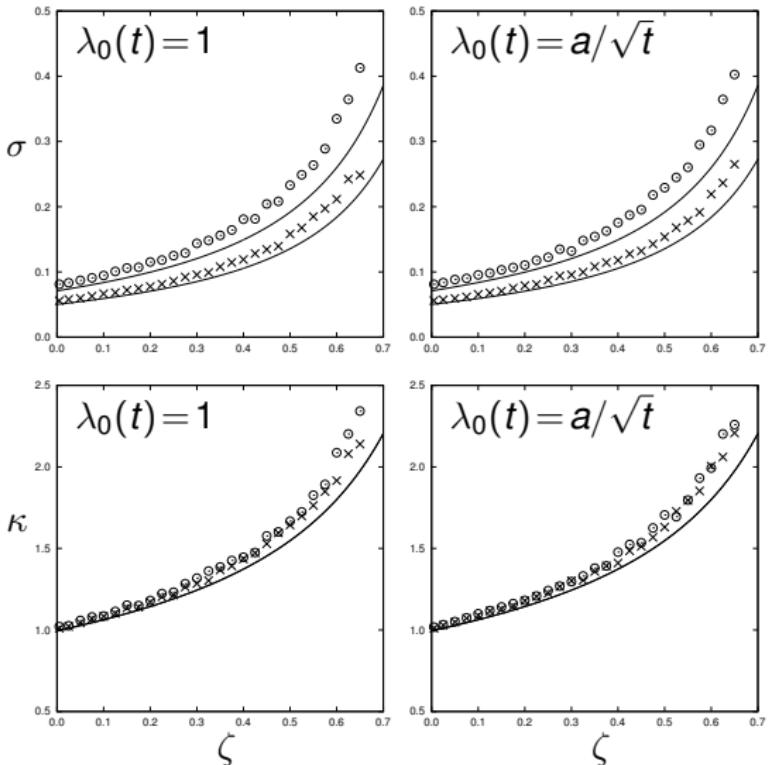
Overfitting in inference with max entropy models

Variational theory: regression parameter clouds

slopes κ and widths σ of regression parameter clouds

$S=0.5$,
 $N=200$ (circles)
or $N=400$ (crosses)

solid lines:
variational theory



Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Summary

- ▶ Overfitting in Cox regression causes predictable bias
 - (i) inflation of association parameters
 - (ii) hazard rates: underestimated (t small), overestimated (t large)
- ▶ Analytical approach based on statistical mechanics (replica method)
 - exact RS equations: $\{u, v, w, \lambda(t)\}$, nontrivial to solve numerically
 - variational approximation: $\{u, v, w, k, \rho\}$, easy to solve numerically
 - predictions of variational theory: quite good,
reliable basis for overfitting corrections
- ▶ Next
 - ▶ Include censoring ✓
 - ▶ Include priors/regularizers ✓
 - ▶ Analysis of exact equations (no variational approx)
 - ▶ Higher order variational approx
 - ▶ Associations for which $\sum_{\mu} \beta_{\mu} z_{\mu}$ is not Gaussian
 - ▶ Roll out overfitting correction protocols for Cox regression ✓
 - ▶ Other survival analysis models

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

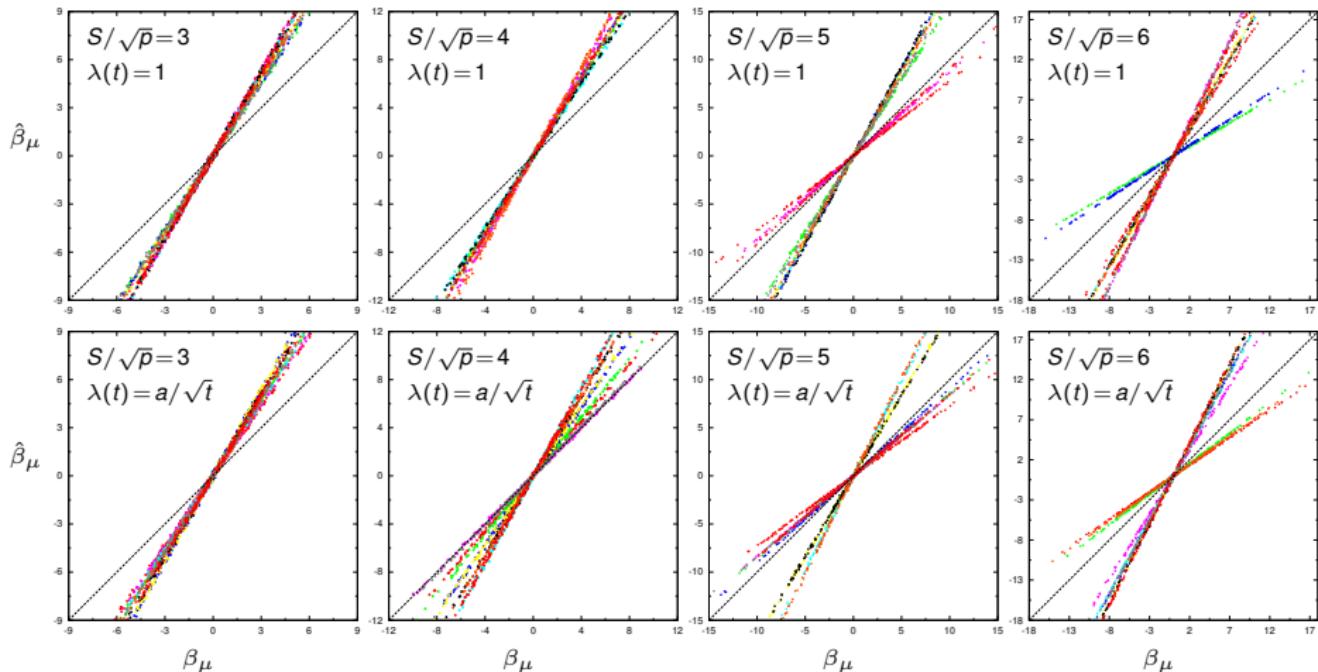
Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Beyond replica symmetry



Large p/N and $\langle \beta_\mu^2 \rangle$: replica symmetry breaking

all cases: $N = 500, \zeta = p/N = 0.4$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Censoring, competing risks and priors

- ▶ inclusion of censoring: essential for application to real data
- ▶ mathematically: censoring and multiple risks very similar
- ▶ inclusion of priors: MAP instead of ML regression

new set-up:

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

$r_i = 1$: *primary risk*
 $r_i = 0$: *end of trial censoring*
 $r_i > 1$: *other active risks*

$$\boldsymbol{\theta}_{\text{MAP}} = \operatorname{argmin}_{\boldsymbol{\theta}} \left[D(\hat{P}_{\mathcal{D}} || P_{\boldsymbol{\theta}}) - \frac{1}{N} \log p(\boldsymbol{\theta}) \right]$$

$$P(\mathcal{D} | \boldsymbol{\theta}) = \prod_{i=1}^N P(t_i, r_i | \mathbf{z}_i, \boldsymbol{\theta})$$

Multiple risk version of Cox model:

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

- ▶ prior: $p(\beta) = Z^{-1} \prod_{\mu} e^{-\phi(\beta_{\mu}/\kappa)}$

theory will now involve covariate correlations,

$$A_{\mu\nu} = \langle z_{\mu} z_{\nu} \rangle - \langle z_{\mu} \rangle \langle z_{\nu} \rangle$$

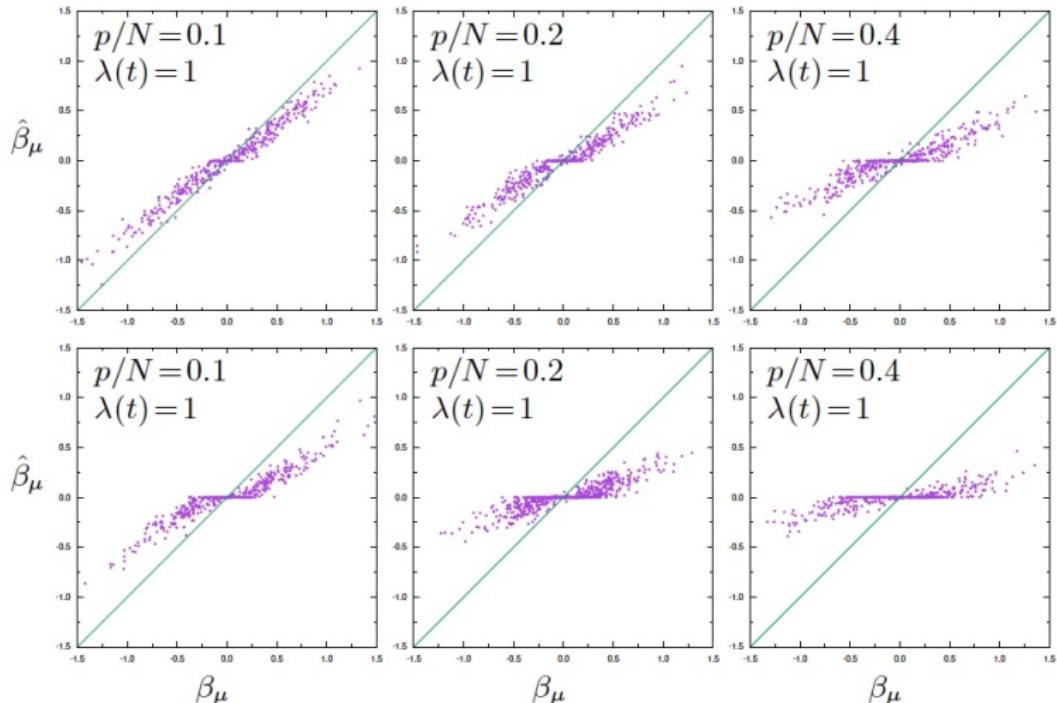
- ▶ more complex stat mech problem,
but still doable:

$$E(\{\beta^0, \lambda^0\}) = \lim_{\gamma \rightarrow \infty} \frac{\partial}{\partial \gamma} \lim_{n \rightarrow 0} \frac{1}{n} \text{extr}_{\Xi, D, d, \lambda_{11}, \dots, \lambda_{Rn}} \Psi[\Xi, D, d, \lambda_{11}, \dots, \lambda_{Rn}]$$

$$\begin{aligned} \Psi[\dots] = & -\frac{1}{2} \zeta \sum_{rs=1}^R \left[\sum_{\alpha\rho=1}^n D_{r\alpha, s\rho} \Xi_{r\alpha, s\rho} + \sum_{\alpha=1}^n \left(d_{rs, \alpha}^a \Xi_{r\alpha, s0} + d_{rs, \alpha}^b \Xi_{r0, s\alpha} \right) \right] \\ & - \lim_{p \rightarrow \infty} \frac{\zeta}{p} \log \int \left[\prod_{r=1}^R \prod_{\alpha=1}^n d\boldsymbol{\eta}^{r\alpha} e^{\gamma \sum_{\mu \leq p} [\phi(\eta_{\mu}^*/\kappa) - \phi(\eta_{\mu}^{r\alpha}/\kappa)]} \right] \\ & \times e^{\sum_{rs=1}^R \left[-i \sum_{\alpha\rho=1}^n \hat{\Xi}_{r\alpha, s\rho} \boldsymbol{\eta}^{r\alpha} \cdot \mathbf{A} \boldsymbol{\eta}^{s\rho} - i \sum_{\alpha=1}^n (\hat{\Xi}_{r\alpha, s0} \boldsymbol{\eta}^{r\alpha} \cdot \mathbf{A} \boldsymbol{\eta}^{s0} + \hat{\Xi}_{r0, s\alpha} \boldsymbol{\eta}^{r0} \cdot \mathbf{A} \boldsymbol{\eta}^{s\alpha}) \right]} \\ - \log \int d\mathbf{Y} \frac{e^{-\frac{1}{2} \sum_{rs=1}^R \sum_{\alpha\rho=0}^n Y_{r\alpha} (\Xi^{-1})_{r\alpha, s\rho} Y_{s\rho}}}{\sqrt{(2\pi)^{(n+1)R} \text{Det } \Xi}} \int dt \sum_{r=0}^R p(t, r | \{y_0, \lambda_0\}) \prod_{\alpha=1}^n \left[\frac{p(t, r | \{y_{\alpha}, \lambda_{\alpha}\})}{p(t, r | \{y_0, \lambda_0\})} \right]^{\gamma} \end{aligned}$$

non-Gaussian priors:
regression parameter clouds no longer linear

e.g. $p(\beta) = Z^{-1} \prod_{\mu} e^{-|\beta_{\mu}|/\sigma}$



simplest case

Gaussian prior, $p(\beta) = (\kappa/2\pi)^{p/2} e^{-\frac{1}{2}\kappa\beta^2}$,
only end-of-trial censoring

- ▶ only spectrum $\varrho(a)$ of covariate correlation matrix
 \mathbf{A} appears in final RS equations
- ▶ explicit expression for $\lambda_0(t)$ of censoring risk
- ▶ two new scalar order parameters
- ▶ RS saddle point problem: extremize

$$\begin{aligned}\Psi_{\text{RS}}[\dots] = & -\frac{1}{2}\zeta \left[g(u^2 + v^2 + w^2) + \gamma\kappa S_0^2 - \langle \log(\gamma\kappa + ag) \rangle \right. \\ & \left. - f\langle \frac{a}{\gamma\kappa + ag} \rangle - \left(\frac{wS_1}{S_0} \right)^2 \langle \frac{a^2}{\gamma\kappa + ag} \rangle^{-1} \right] \\ & + \gamma \int Dy_0 \int dt \sum_{r=0}^1 p(t, r | S_1 y_0, \lambda^*) \log p(t, r | S_1 y_0, \lambda^*) \\ & - \int Dy_0 Dz \int dt \sum_{r=0}^1 p(t, r | S_1 y_0, \lambda^*) \log \int Dy p^\gamma(t, r | uy + vz + wy_0, \lambda)\end{aligned}$$

with: $S_\ell^2 = \frac{1}{p} \beta^* \mathbf{A}^\ell \beta^*$, $\langle g(a) \rangle = \int da \varrho(a) g(a)$

Introduction

Overfitting in Cox regression

Simulations

Quantifying overfitting in Cox regression

Intuition for the problem

The basic ideas

Translation to Cox's model

Statmech analysis

Conversion to saddle-point problem

Replica symmetric extrema

Physical interpretation of order parameters

RS saddle point equations

Analysis of RS equations

RS eqns in the limit $\gamma \rightarrow \infty$

Numerical and asymptotic solution of RS eqns

Variational approximation

Tests and applications

Variational theory: tests of order parameters

Variational theory: regression parameter clouds

Summary and ongoing work

Summary

Beyond replica symmetry

Censoring, competing risks and priors

Overfitting in inference with max entropy models

Overfitting in inference with max entropy models

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

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

- ▶ Exponential models:

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

e.g: inference of

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

- ▶ MAP inference, with Gaussian priors:

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

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