

## Focused Regularised Likelihood

Gudmund Horn Hermansen with Nils Lid Hjort

# Godt Hjort

Næsten alle er enige om at det er godt, men den med 4  
med mange ulige arter godt, men den med 4  
dårlig.

( fjerner man denne passer dit veldig fint, noe som  
kan være med på å gi modellen troverdighet )

Mye bra her, Gudmund. Ta hovedfag.  
Nils

OK diskusjon.

## Introduction, Summary and Notation

Suppose we have data from

$$Y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_p x_{i,p} + \sigma \epsilon_i = x_i^t \beta + \sigma \epsilon_i,$$

with  $\epsilon_1, \dots, \epsilon_n$  are i.i.d. standard normals,  $\beta \in \mathbb{R}^{p+1}$  and  $\sigma > 0$ .

## Introduction, Summary and Notation

Suppose we have data from

$$Y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_p x_{i,p} + \sigma \epsilon_i = x_i^t \beta + \sigma \epsilon_i,$$

with  $\epsilon_1, \dots, \epsilon_n$  are i.i.d. standard normals,  $\beta \in \mathbb{R}^{p+1}$  and  $\sigma > 0$ .

**Ridge regression** is a common regularised method for estimating  $\beta$ :

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - x_i^t \beta)^2 + \lambda \sum_{j=1}^{p+1} |\beta_j|^2 \right\}.$$

Note that:

- $\lambda = 0$  is the same as **ordinary least squares regression**
- increasing  $\lambda$  will ‘shrink’ the  $\hat{\beta}_j$ -s toward zero

## Introduction, Summary and Notation

Suppose we have data from

$$Y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_p x_{i,p} + \sigma \epsilon_i = x_i^t \beta + \sigma \epsilon_i,$$

with  $\epsilon_1, \dots, \epsilon_n$  are i.i.d. standard normals,  $\beta \in \mathbb{R}^{p+1}$  and  $\sigma > 0$ .

**Ridge regression** is a common regularised method for estimating  $\beta$ :

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - x_i^t \beta)^2 + \lambda \sum_{j=1}^{p+1} |\beta_j|^2 \right\}.$$

Note that:

- $\lambda = 0$  is the same as **ordinary least squares regression**
- increasing  $\lambda$  will 'shrink' the  $\hat{\beta}_j$ -s toward zero

The above is (essentially) the same as

$$\hat{\theta}_{\lambda} = (\hat{\beta}_{\lambda}, \hat{\sigma}_{\lambda}) = \arg \max_{\beta, \sigma} \left\{ \ell_n(\beta, \sigma) + \lambda n \sum_{j=1}^{p+1} |\beta_j|^2 \right\},$$

where  $\ell_n(\beta)$  is the log-likelihood corresponding to a Gaussian distribution.

**The penalisation term is now scaled by  $n$ .**

## Introduction, Summary and Notation

We denote the ‘true’ data-generating distribution by  $G$ .

In general  $Y_1, Y_2, \dots, Y_n$  are i.i.d. from  $G$ .

Let  $F_\theta$  be a parametric model, and  $\ell_n(\theta)$  the corresponding log-likelihood.

## Introduction, Summary and Notation

We denote the ‘true’ data-generating distribution by  $G$ .

In general  $Y_1, Y_2, \dots, Y_n$  are i.i.d. from  $G$ .

Let  $F_\theta$  be a parametric model, and  $\ell_n(\theta)$  the corresponding log-likelihood.

Then the **FRL estimator** is

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\},$$

where  $\lambda$  is a tuning parameter and  $\psi$  a **control parameter**.

## Introduction, Summary and Notation

We denote the ‘true’ data-generating distribution by  $G$ .

In general  $Y_1, Y_2, \dots, Y_n$  are i.i.d. from  $G$ .

Let  $F_\theta$  be a parametric model, and  $\ell_n(\theta)$  the corresponding log-likelihood.

Then the **FRL estimator** is

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\},$$

where  $\lambda$  is a tuning parameter and  $\psi$  a **control parameter**.

Effective control parameters are **important characteristics** of a distribution where we also have **robust alternative estimators** (non-parametric).

**Example:** A quantile with  $0 \leq p \leq 1$ :

$$\psi(\theta) = \psi(F_\theta, p) = F_\theta^{-1}(p) \quad \text{and} \quad \hat{\psi} = \hat{\psi}(p) = \hat{G}_n^{-1}(p),$$

where  $\hat{G}_n$  is the empirical CDF.



## Introduction, Summary and Notation

We denote the ‘true’ data-generating distribution by  $G$ .

In general  $Y_1, Y_2, \dots, Y_n$  are i.i.d. from  $G$ .

Let  $F_\theta$  be a parametric model, and  $\ell_n(\theta)$  the corresponding log-likelihood.

Then the **FRL estimator** is

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\},$$

where  $\lambda$  is a tuning parameter and  $\psi$  a **control parameter**.

Effective control parameters are **important characteristics** of a distribution where we also have **robust alternative estimators** (non-parametric).

**Example:**  $k$ -th moment:

$$\psi(\theta) = \psi(F_\theta, k) = \int y^k dF_\theta(y) \quad \text{and} \quad \hat{\psi} = \hat{\psi}(k) = \frac{1}{n} \sum_{i=1}^n y_i^k.$$

## Introduction, Summary and Notation

We denote the ‘true’ data-generating distribution by  $G$ .

In general  $Y_1, Y_2, \dots, Y_n$  are i.i.d. from  $G$ .

Let  $F_\theta$  be a parametric model, and  $\ell_n(\theta)$  the corresponding log-likelihood.

Then the **FRL estimator** is

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\},$$

where  $\lambda$  is a tuning parameter and  $\psi$  a **control parameter**.

Effective control parameters are **important characteristics** of a distribution where we also have **robust alternative estimators** (non-parametric).

**Example:** A probability, e.g.

$$\psi(\theta) = \psi(F_\theta, q) = \int I(y > q) dF_\theta(y) \quad \text{and} \quad \hat{\psi}(q) = \frac{1}{n} \sum_{i=1}^n I(y_i > q).$$

In general, the log Focused Regularised Likelihood (log-FRL) is

$$\ell_{n,\lambda}(\theta) = \ell_{n,\lambda,\mathbf{w},\boldsymbol{\psi}}(\theta) = \ell_n(\theta) - \frac{1}{2}\lambda n \sum_{j=1}^r w_j \{\widehat{\psi}_j - \psi_j(\theta)\}^2,$$

where

- $\ell_n(\theta)$  is the log-likelihood corresponding to  $F_\theta$
- $\lambda$  is a tuning parameter
- $\psi_j$  are control or focus parameter, e.g. quantiles, moments, ...
- $\widehat{\psi}_j$  are non-parametric or robust alternative estimates for  $\psi_j$
- $w_1, \dots, w_r$  are weights with  $w_1 + \dots + w_r = 1$

## Introduction, Summary and Notation

In general, the log Focused Regularised Likelihood (log-FRL) is

$$\ell_{n,\lambda}(\theta) = \ell_{n,\lambda,w,\psi}(\theta) = \ell_n(\theta) - \frac{1}{2}\lambda n \sum_{j=1}^r w_j \{\widehat{\psi}_j - \psi_j(\theta)\}^2,$$

where

- $\ell_n(\theta)$  is the log-likelihood corresponding to  $F_\theta$
- $\lambda$  is a tuning parameter
- $\psi_j$  are control or focus parameter, e.g. quantiles, moments, ...
- $\widehat{\psi}_j$  are non-parametric or robust alternative estimates for  $\psi_j$
- $w_1, \dots, w_r$  are weights with  $w_1 + \dots + w_r = 1$

Note that:

- if  $\lambda = 0$  we have  $\widehat{\theta}_\lambda = \widehat{\theta}_{\text{ML}}$
- increasing  $\lambda$  will 'push'  $\psi_j(\widehat{\theta}_\lambda)$  to match  $\widehat{\psi}_j$

We also need a set of 'standard' **regularity assumptions** to be true.

## Why?

Control parameters can make the estimated model **more robust**.

Control as a **focus** parameter can improve the model where it is **important**.

# Why?

Control parameters can make the estimated model **more robust**.

Control as a **focus** parameter can improve the model where it is **important**.

Analytic large sample theory for:

- standard models for i.i.d. data
- **models with local misspecification**
- regression models
- stationary time series (will not talk about this here)

## Illustration: Measuring the Speed of Light

We can use FRL for **robust estimation** of a normal density.

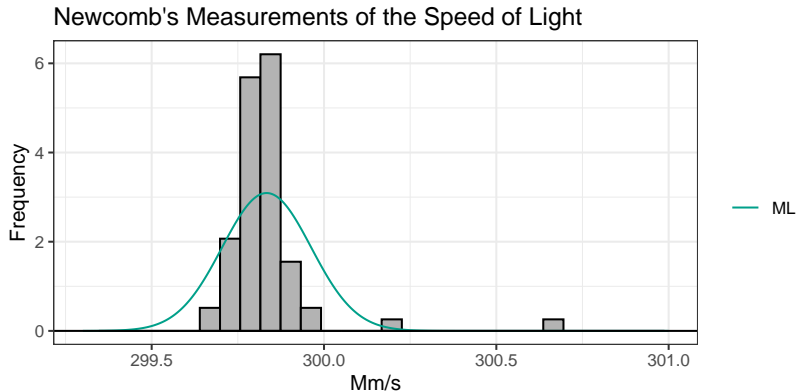
In particular if data contains **outliers or is contaminated**.

## Illustration: Measuring the Speed of Light

We can use FRL for **robust estimation** of a normal density.

In particular if data contains **outliers or is contaminated**.

Suppose we want to model the data below with a  $N(\mu, \sigma^2)$ .



Simon Newcomb speed of light measurements; see e.g. Stigler (1977) for details about the data.



## Illustration: Measuring the Speed of Light

We will do this by adding some **quantiles as control parameters**:

$$\psi(\mu, \sigma, p) = \sigma \times \Phi^{-1}(p) + \mu \quad \text{and} \quad \hat{\psi} = \hat{\psi}(p) = \hat{G}^{-1}(p)$$

where  $\hat{G}$  is the empirical CDF, with  $p = 0.1, 0.5, p = 0.90$  and  $\lambda = 1000$ .

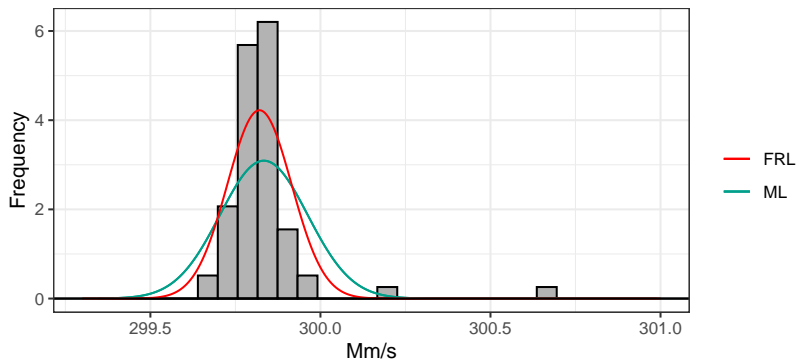
## Illustration: Measuring the Speed of Light

We will do this by adding some **quantiles as control parameters**:

$$\psi(\mu, \sigma, p) = \sigma \times \Phi^{-1}(p) + \mu \quad \text{and} \quad \hat{\psi} = \hat{\psi}(p) = \hat{G}^{-1}(p)$$

where  $\hat{G}$  is the empirical CDF, with  $p = 0.1, 0.5, p = 0.90$  and  $\lambda = 1000$ .

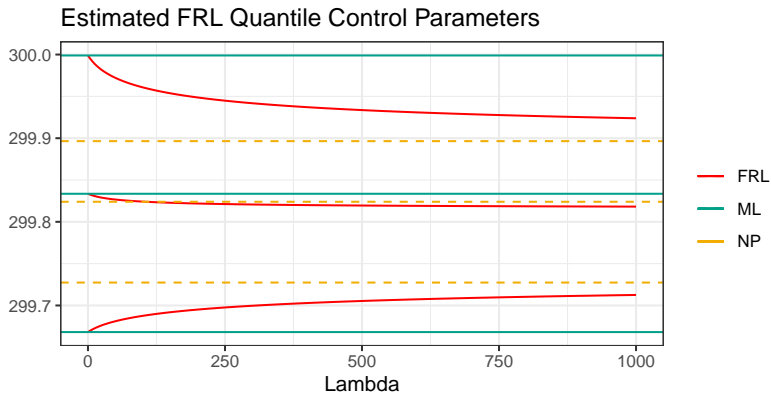
Newcomb's Measurements of the Speed of Light



Simon Newcomb speed of light measurements; see e.g. Stigler (1977) for details about the data.

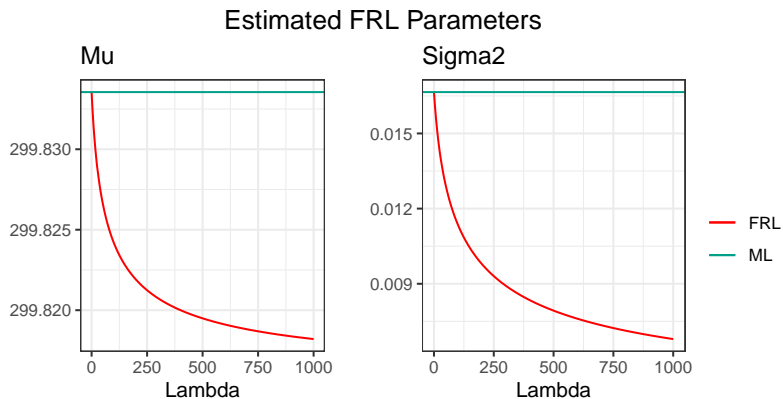
## Illustration: Measuring the Speed of Light

Increasing  $\lambda$  'push' the estimated quantiles towards **the empirical quantiles**.



## Illustration: Measuring the Speed of Light

And, the estimated parameters move away from the ML estimates.



However, are these FRL estimates more precise?

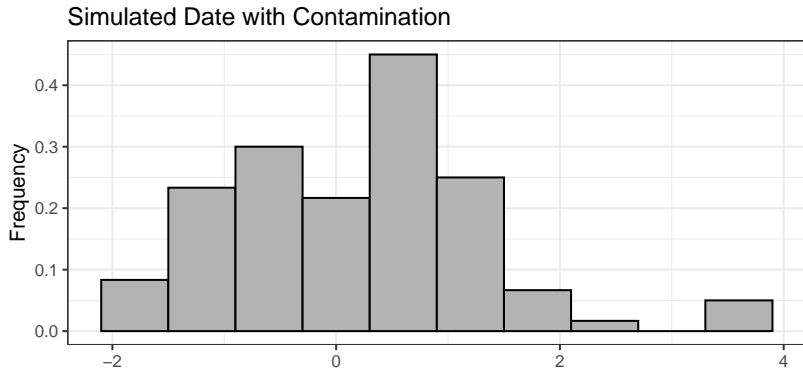
## Illustration: Simulated Data with Contamination

Simulated data **with contamination** (outliers).

Repeated simulations of independent  $Y_1, \dots, Y_{100}$  with  $Y_i \sim N(0, 1)$ .

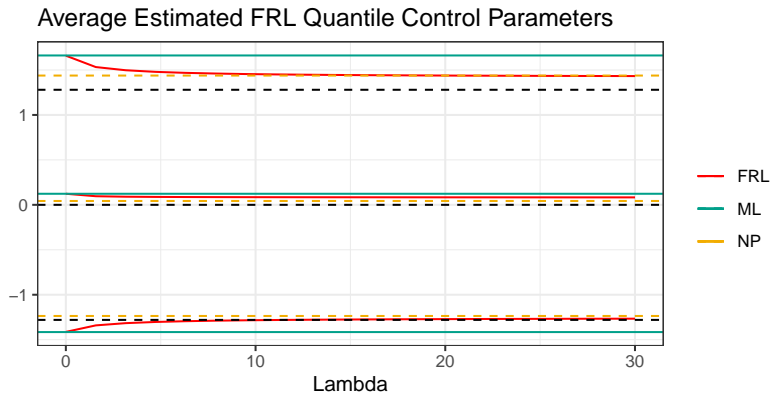
Add 4% contamination from a  $N(4, 0.5)$ .

Again, we will use **control parameters based on quantiles** (0.1, 0.5 and 0.9).



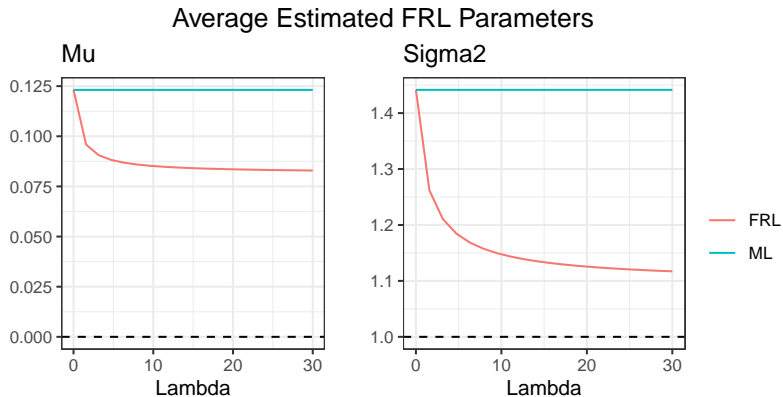
## Illustration: Simulated Data with Contamination

Estimated quantiles are 'pushed' towards the empirical (and true) quantiles.



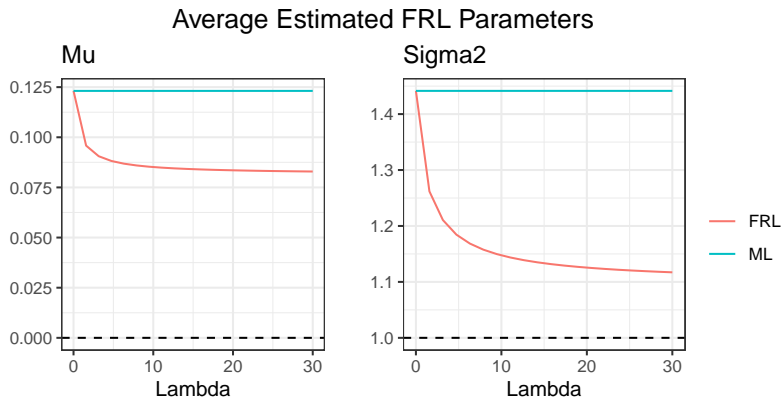
## Illustration: Simulated Data with Contamination

And the estimated parameters move closer to **the true values**.



## Illustration: Simulated Data with Contamination

And the estimated parameters move closer to **the true values**.



Here, the **median** was used as a control parameter.

However, should we just **use the non-parametric estimate(s)**?



## Illustration: Estimation of Location Parameter with Contamination

Same simulation setup, but with ‘focus’ on estimating a **location parameter**.

## Illustration: Estimation of Location Parameter with Contamination

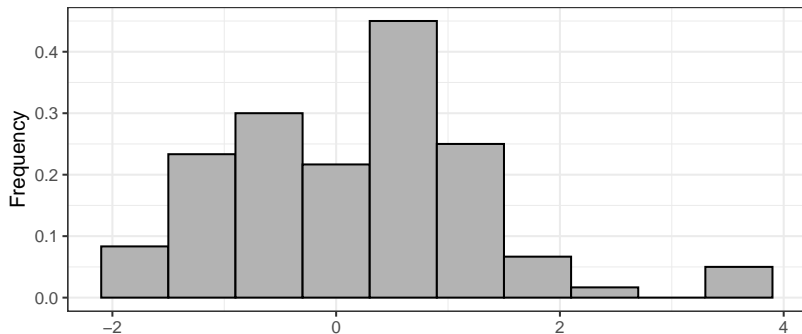
Same simulation setup, but with ‘focus’ on estimating a **location parameter**.

We frame this as estimating  $\mu$  in a  $N(\mu, 1)$ , two natural estimators are

$$\hat{\mu}_{\text{ML}} = \bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{and} \quad \hat{m} = \text{median}(Y_1, \dots, Y_n).$$

We consider  $\mu = 0$  to be the ‘true’ target value.

Simulated Data with Contamination



## Illustration: Estimation of Location Parameter with Contamination

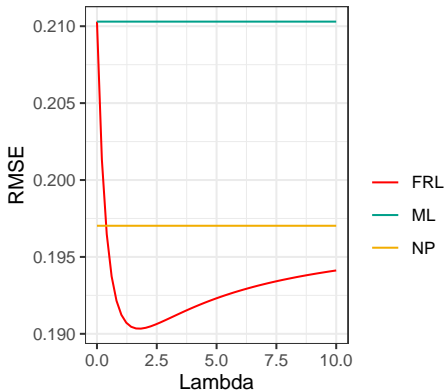
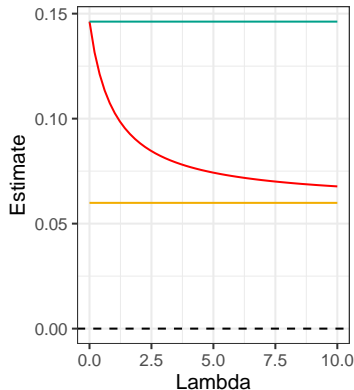
The FRL, with the **median as control**  $\psi = \mu$  and  $\hat{\psi} = \hat{m}$  is

$$\hat{\mu}_\lambda = \arg \max_{\mu} \left\{ \ell_n(\mu) - \frac{1}{2} \lambda n \{ \hat{m} - \mu \}^2 \right\}$$

## Illustration: Estimation of Location Parameter with Contamination

The FRL, with the **median as control**  $\psi = \mu$  and  $\hat{\psi} = \hat{m}$  is

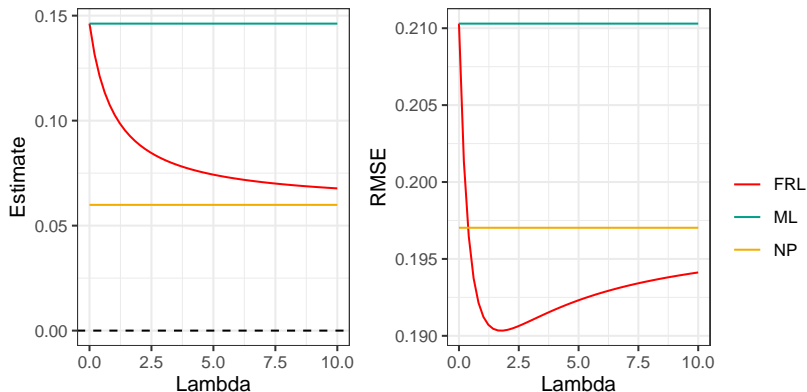
$$\hat{\mu}_\lambda = \arg \max_{\mu} \left\{ \ell_n(\mu) - \frac{1}{2} \lambda n \{ \hat{m} - \mu \}^2 \right\}$$



## Illustration: Estimation of Location Parameter with Contamination

The FRL, with the **median as control**  $\psi = \mu$  and  $\hat{\psi} = \hat{m}$  is

$$\hat{\mu}_\lambda = \arg \max_{\mu} \left\{ \ell_n(\mu) - \frac{1}{2} \lambda n \{ \hat{m} - \mu \}^2 \right\}$$



Note that this is a **bias–variance trade-off** game (with respect to RMSE).

How to determine the **optimal**  $\lambda$ ?

## Large Sample Theory - Illustration

We can use **large-sample theory** or the **bootstrap** to analyse the ‘behaviour’ of the FRL estimator, find optimal  $\lambda$ , compare it to the MLE, . . . .

A bootstrap approach seems to work well, will not focus on this here.

## Large Sample Theory - Illustration

We can use **large-sample theory** or the **bootstrap** to analyse the ‘behaviour’ of the FRL estimator, find optimal  $\lambda$ , compare it to the MLE, . . . .

A bootstrap approach seems to work well, will not focus on this here.

We need **the target of  $\hat{\theta}_\lambda$** , say  $\theta_\lambda$ , and the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

## Large Sample Theory - Illustration

We can use **large-sample theory** or the **bootstrap** to analyse the ‘behaviour’ of the FRL estimator, find optimal  $\lambda$ , compare it to the MLE, ...

A bootstrap approach seems to work well, will not focus on this here.

We need **the target of  $\hat{\theta}_\lambda$** , say  $\theta_\lambda$ , and the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

Again, consider estimating a **location parameter**, with competing estimators

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{and} \quad \hat{m} = \text{median}(Y_1, \dots, Y_n),$$

with  $Y_i$  are i.i.d. and  $Y_i \sim G$ .

In order to fit the FRL framework, we view this as estimating  $\mu$  in a  $N(\mu, 1)$ .



## Large Sample Theory - Illustration

We can use **large-sample theory** or the **bootstrap** to analyse the ‘behaviour’ of the FRL estimator, find optimal  $\lambda$ , compare it to the MLE, ...

A bootstrap approach seems to work well, will not focus on this here.

We need **the target of  $\hat{\theta}_\lambda$** , say  $\theta_\lambda$ , and the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

Again, consider estimating a **location parameter**, with competing estimators

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{and} \quad \hat{m} = \text{median}(Y_1, \dots, Y_n),$$

with  $Y_i$  are i.i.d. and  $Y_i \sim G$ .

In order to fit the FRL framework, we view this as estimating  $\mu$  in a  $N(\mu, 1)$ .

With the **median as the control parameter**, i.e.  $\psi(\mu) = \mu$  and  $\hat{\psi} = \hat{m}$  as the robust alternative, then

$$\hat{\mu}_\lambda = \arg \max_{\mu} \left\{ \ell_n(\mu) - \frac{1}{2} \lambda n \{ \hat{m} - \mu \}^2 \right\}$$

## Large Sample Theory - Illustration

We can use **large-sample theory** or the **bootstrap** to analyse the ‘behaviour’ of the FRL estimator, find optimal  $\lambda$ , compare it to the MLE, ...

A bootstrap approach seems to work well, will not focus on this here.

We need **the target of  $\hat{\theta}_\lambda$** , say  $\theta_\lambda$ , and the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

Again, consider estimating a **location parameter**, with competing estimators

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{and} \quad \hat{m} = \text{median}(Y_1, \dots, Y_n),$$

with  $Y_i$  are i.i.d. and  $Y_i \sim G$ .

In order to fit the FRL framework, we view this as estimating  $\mu$  in a  $N(\mu, 1)$ .

With the **median as the control parameter**, i.e.  $\psi(\mu) = \mu$  and  $\hat{\psi} = \hat{m}$  as the robust alternative, then

$$\hat{\mu}_\lambda = \arg \max_{\mu} \left\{ \ell_n(\mu) - \frac{1}{2} \lambda n \{ \hat{m} - \mu \}^2 \right\}$$

and

$$\hat{\mu}_\lambda = \frac{1}{1 + \lambda} \bar{Y}_n + \frac{\lambda}{1 + \lambda} \hat{m}.$$

## Large Sample Theory - Illustration

We know that  $\bar{Y}_n \rightarrow_{\text{pr}} E_g Y_1 = \mu_g$ , and  $\hat{m} \rightarrow_{\text{pr}} m_g = G^{-1}(0.5)$ .

## Large Sample Theory - Illustration

We know that  $\bar{Y}_n \rightarrow_{\text{pr}} E_g Y_1 = \mu_g$ , and  $\hat{m} \rightarrow_{\text{pr}} m_g = G^{-1}(0.5)$ .

And, we can show that

$$\hat{\mu} \rightarrow_{\text{pr}} \mu_\lambda = \arg \min_{\mu} \left\{ \text{KL}(g, f_{\mu}) + \frac{1}{2} \lambda \{m - \mu\}^2 \right\} = \frac{1}{1 + \lambda} \mu_g + \frac{\lambda}{1 + \lambda} m_g,$$

where  $f_{\mu}$  is the density of a  $N(\mu, 1)$ .

## Large Sample Theory - Illustration

We know that  $\bar{Y}_n \rightarrow_{\text{pr}} E_g Y_1 = \mu_g$ , and  $\hat{m} \rightarrow_{\text{pr}} m_g = G^{-1}(0.5)$ .

And, we can show that

$$\hat{\mu} \rightarrow_{\text{pr}} \mu_\lambda = \arg \min_{\mu} \left\{ \text{KL}(g, f_{\mu}) + \frac{1}{2} \lambda \{m - \mu\}^2 \right\} = \frac{1}{1 + \lambda} \mu_g + \frac{\lambda}{1 + \lambda} m_g,$$

where  $f_{\mu}$  is the density of a  $N(\mu, 1)$ .

Moreover,

$$\sqrt{n}(\hat{\mu} - \mu_\lambda) \rightarrow_d \Lambda \sim N \left( 0, \frac{1}{(1 + \lambda)^2} \left[ \sigma_g^2 + \frac{\lambda^2}{4g(m_g)^2} + \frac{\lambda \times E_g |Y_1 - m_g|}{g(m_g)} \right] \right),$$

where  $\sigma_g^2 = \text{Var}_g(Y_1)$ .

## Large Sample Theory - Illustration

We know that  $\bar{Y}_n \rightarrow_{\text{pr}} E_g Y_1 = \mu_g$ , and  $\hat{m} \rightarrow_{\text{pr}} m_g = G^{-1}(0.5)$ .

And, we can show that

$$\hat{\mu} \rightarrow_{\text{pr}} \mu_\lambda = \arg \min_{\mu} \left\{ \text{KL}(g, f_{\mu}) + \frac{1}{2} \lambda \{m - \mu\}^2 \right\} = \frac{1}{1 + \lambda} \mu_g + \frac{\lambda}{1 + \lambda} m_g,$$

where  $f_{\mu}$  is the density of a  $N(\mu, 1)$ .

Moreover,

$$\sqrt{n}(\hat{\mu} - \mu_\lambda) \rightarrow_d \Lambda \sim N \left( 0, \frac{1}{(1 + \lambda)^2} \left[ \sigma_g^2 + \frac{\lambda^2}{4g(m_g)^2} + \frac{\lambda \times E_g |Y_1 - m_g|}{g(m_g)} \right] \right),$$

where  $\sigma_g^2 = \text{Var}_g(Y_1)$ .

There are **analogous limit distribution** results for both  $\bar{Y}$  and the median  $\hat{m}$ .

From this, we can extract (limit) **bias, variance, RMSE**, etc.

## Large Sample Theory - Illustration

We know that  $\bar{Y}_n \rightarrow_{\text{pr}} E_g Y_1 = \mu_g$ , and  $\hat{m} \rightarrow_{\text{pr}} m_g = G^{-1}(0.5)$ .

And, we can show that

$$\hat{\mu} \rightarrow_{\text{pr}} \mu_\lambda = \arg \min_{\mu} \left\{ \text{KL}(g, f_{\mu}) + \frac{1}{2} \lambda \{m - \mu\}^2 \right\} = \frac{1}{1 + \lambda} \mu_g + \frac{\lambda}{1 + \lambda} m_g,$$

where  $f_{\mu}$  is the density of a  $N(\mu, 1)$ .

Moreover,

$$\sqrt{n}(\hat{\mu} - \mu_\lambda) \rightarrow_d \Lambda \sim N \left( 0, \frac{1}{(1 + \lambda)^2} \left[ \sigma_g^2 + \frac{\lambda^2}{4g(m_g)^2} + \frac{\lambda \times E_g |Y_1 - m_g|}{g(m_g)} \right] \right),$$

where  $\sigma_g^2 = \text{Var}_g(Y_1)$ .

There are **analogous limit distribution** results for both  $\bar{Y}$  and the median  $\hat{m}$ .

From this, we can extract (limit) **bias, variance, RMSE**, etc.

And, for example:

- we can derive an expression for the optimal value of  $\lambda$
- compare estimators
- make asymptotic test and diagnostics tools/plots

## Large Sample Theory - Summary

Consider one  $\psi$  and let  $Y_1, Y_2, \dots, Y_n$  be i.i.d. from  $G$ , then

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\}.$$

In order to ‘understand’ the FRL estimate we need to:

- (1) find what  $\hat{\theta}_\lambda$  **aims at** and
- (2) derive the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .



## Large Sample Theory - Summary

Consider one  $\psi$  and let  $Y_1, Y_2, \dots, Y_n$  be i.i.d. from  $G$ , then

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\}.$$

In order to ‘understand’ the FRL estimate we need to:

- (1) find what  $\hat{\theta}_\lambda$  **aims at** and
- (2) derive the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

We obtain (1) by similar arguments as the MLE for a misspecified model:

$$\hat{\theta}_\lambda \rightarrow_{\text{pr}} \theta_\lambda = \arg \min_{\theta} \left\{ \text{KL}(g, f_\theta) + \frac{1}{2} \lambda \{ \psi_{\text{true}} - \psi(\theta) \}^2 \right\}.$$

## Large Sample Theory - Summary

Consider one  $\psi$  and let  $Y_1, Y_2, \dots, Y_n$  be i.i.d. from  $G$ , then

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\}.$$

In order to ‘understand’ the FRL estimate we need to:

- (1) find what  $\hat{\theta}_\lambda$  **aims at** and
- (2) derive the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

We obtain (1) by similar arguments as the MLE for a misspecified model:

$$\hat{\theta}_\lambda \rightarrow_{\text{pr}} \theta_\lambda = \arg \min_{\theta} \left\{ \text{KL}(g, f_\theta) + \frac{1}{2} \lambda \{ \psi_{\text{true}} - \psi(\theta) \}^2 \right\}.$$

Similar for (2), where we can show

$$\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda) \rightarrow_d N_p(0, [J(\theta_\lambda) + \lambda L]^{-1} K_\lambda [J(\theta_\lambda) + \lambda L]^{-1}),$$

## Large Sample Theory - Summary

Consider one  $\psi$  and let  $Y_1, Y_2, \dots, Y_n$  be i.i.d. from  $G$ , then

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi} - \psi(\theta) \}^2 \right\}.$$

In order to ‘understand’ the FRL estimate we need to:

- (1) find what  $\hat{\theta}_\lambda$  **aims at** and
- (2) derive the **limit distribution** of  $\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda)$ .

We obtain (1) by similar arguments as the MLE for a misspecified model:

$$\hat{\theta}_\lambda \rightarrow_{\text{pr}} \theta_\lambda = \arg \min_{\theta} \left\{ \text{KL}(g, f_{\theta}) + \frac{1}{2} \lambda \{ \psi_{\text{true}} - \psi(\theta) \}^2 \right\}.$$

Similar for (2), where we can show

$$\sqrt{n}(\hat{\theta}_\lambda - \theta_\lambda) \rightarrow_d N_p(0, [J(\theta_\lambda) + \lambda L]^{-1} K_{\lambda} [J(\theta_\lambda) + \lambda L]^{-1}),$$

where  $J$  is the Fisher information,  $L = \dot{\psi}(\theta_\lambda) \dot{\psi}(\theta_\lambda)^t + [\psi_{\text{true}} - \psi(\theta_\lambda)] \ddot{\psi}(\theta_\lambda)$ ,

$$K_{\lambda} = K(\theta_\lambda) + 2\lambda c \dot{\psi}(\theta_\lambda)^t + \lambda^2 \tau^2 \dot{\psi}(\theta_\lambda) \dot{\psi}(\theta_\lambda)^t$$

and  $K(\cdot)$ ,  $\tau^2$  and  $c$  are elements from the covariance matrix involving  $\hat{\psi}$  and the scaled score-function,  $\dot{\psi}$  and  $\ddot{\psi}$  are first and second order derivatives.

## Models with Local Misspecification - Summary

Suppose  $Y_1, Y_2, \dots, Y_n$  are i.i.d. and  $Y_i \sim F_{\theta_0, \gamma_0 + \delta/\sqrt{n}}$  and  $\gamma_0$  is known.

## Models with Local Misspecification - Summary

Suppose  $Y_1, Y_2, \dots, Y_n$  are i.i.d. and  $Y_i \sim F_{\theta_0, \gamma_0 + \delta/\sqrt{n}}$  and  $\gamma_0$  is known.

Models with local misspecification are useful for examining bias–variance trade–offs in a large-sample framework; Claeskens and Hjort (2008).

## Models with Local Misspecification - Summary

Suppose  $Y_1, Y_2, \dots, Y_n$  are i.i.d. and  $Y_i \sim F_{\theta_0, \gamma_0 + \delta/\sqrt{n}}$  and  $\gamma_0$  is known.

Models with local misspecification are useful for examining bias–variance trade–offs in a large-sample framework; Claeskens and Hjort (2008).

Assume we are interested in estimating  $\psi = \psi(\theta, \gamma)$  (focus parameter).

## Models with Local Misspecification - Summary

Suppose  $Y_1, Y_2, \dots, Y_n$  are i.i.d. and  $Y_i \sim F_{\theta_0, \gamma_0 + \delta/\sqrt{n}}$  and  $\gamma_0$  is known.

**Models with local misspecification** are useful for examining **bias–variance trade-offs** in a large-sample framework; Claeskens and Hjort (2008).

Assume we are interested in estimating  $\psi = \psi(\theta, \gamma)$  (**focus parameter**).

We can compare  $\psi_{\text{narr}} = \psi(\hat{\theta}_{\text{narr}}, \gamma_0)$  with  $\psi_{\text{wide}} = \psi(\hat{\theta}, \hat{\gamma})$  in the limit, i.e.

$$\sqrt{n}(\hat{\psi}_{\text{narr}} - \psi_{\text{true}}) \rightarrow_d N(\omega\delta, \tau_0^2)$$

$$\sqrt{n}(\hat{\psi}_{\text{wide}} - \psi_{\text{true}}) \rightarrow_d N(0, \tau_0^2 + \omega^2 \kappa^2)$$

with  $\omega = J_{10} J_{00}^{-1} \dot{\psi}_{\theta} - \dot{\psi}_{\gamma}$  and  $\tau_0^2 = \dot{\psi}_{\theta}^t J_{00}^{-1} \dot{\psi}_{\theta}$ , and  $\dot{\psi}$  are partial derivatives.

## Models with Local Misspecification - Summary

Suppose  $Y_1, Y_2, \dots, Y_n$  are i.i.d. and  $Y_i \sim F_{\theta_0, \gamma_0 + \delta/\sqrt{n}}$  and  $\gamma_0$  is known.

**Models with local misspecification** are useful for examining **bias–variance trade-offs** in a large-sample framework; Claeskens and Hjort (2008).

Assume we are interested in estimating  $\psi = \psi(\theta, \gamma)$  (**focus parameter**).

We can compare  $\psi_{\text{narr}} = \psi(\hat{\theta}_{\text{narr}}, \gamma_0)$  with  $\psi_{\text{wide}} = \psi(\hat{\theta}, \hat{\gamma})$  in the limit, i.e.

$$\sqrt{n}(\hat{\psi}_{\text{narr}} - \psi_{\text{true}}) \rightarrow_d N(\omega\delta, \tau_0^2)$$

$$\sqrt{n}(\hat{\psi}_{\text{wide}} - \psi_{\text{true}}) \rightarrow_d N(0, \tau_0^2 + \omega^2 \kappa^2)$$

with  $\omega = J_{10} J_{00}^{-1} \dot{\psi}_{\theta} - \dot{\psi}_{\gamma}$  and  $\tau_0^2 = \dot{\psi}_{\theta}^t J_{00}^{-1} \dot{\psi}_{\theta}$ , and  $\dot{\psi}$  are partial derivatives.

If the **FRL estimate** is

$$\hat{\theta}_{\lambda} = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi}_{\text{wide}} - \psi(\theta, \gamma_0) \}^2 \right\},$$

and  $\hat{\psi}_{\lambda} = \psi(\hat{\theta}_{\lambda}, \gamma_0)$



## Models with Local Misspecification - Summary

Suppose  $Y_1, Y_2, \dots, Y_n$  are i.i.d. and  $Y_i \sim F_{\theta_0, \gamma_0 + \delta/\sqrt{n}}$  and  $\gamma_0$  is known.

**Models with local misspecification** are useful for examining **bias-variance trade-offs** in a large-sample framework; Claeskens and Hjort (2008).

Assume we are interested in estimating  $\psi = \psi(\theta, \gamma)$  (**focus parameter**).

We can compare  $\psi_{\text{narr}} = \psi(\hat{\theta}_{\text{narr}}, \gamma_0)$  with  $\psi_{\text{wide}} = \psi(\hat{\theta}, \hat{\gamma})$  in the limit, i.e.

$$\sqrt{n}(\hat{\psi}_{\text{narr}} - \psi_{\text{true}}) \rightarrow_d N(\omega\delta, \tau_0^2)$$

$$\sqrt{n}(\hat{\psi}_{\text{wide}} - \psi_{\text{true}}) \rightarrow_d N(0, \tau_0^2 + \omega^2 \kappa^2)$$

with  $\omega = J_{10}J_{00}^{-1}\dot{\psi}_\theta - \dot{\psi}_\gamma$  and  $\tau_0^2 = \dot{\psi}_\theta^t J_{00}^{-1} \dot{\psi}_\theta$ , and  $\dot{\psi}$  are partial derivatives.

If the **FRL estimate** is

$$\hat{\theta}_\lambda = \arg \max_{\theta} \left\{ \ell_n(\theta) - \frac{1}{2} \lambda n \{ \hat{\psi}_{\text{wide}} - \psi(\theta, \gamma_0) \}^2 \right\},$$

and  $\hat{\psi}_\lambda = \psi(\hat{\theta}_\lambda, \gamma_0)$  we can show that

$$\sqrt{n}(\hat{\psi}_\lambda - \psi_{\text{true}}) \rightarrow_d N(\omega_\lambda \delta, \tau_\lambda^2),$$

with  $J_\lambda = J_{00} + \lambda \dot{\psi}_\theta \dot{\psi}_\theta^t$  and

$\omega_\lambda = (J_{01} + \lambda \dot{\psi}_\gamma \dot{\psi}_\theta) J_\lambda^{-1} \dot{\psi}_\theta - \dot{\psi}_\gamma$  and  $\tau_\lambda^2 = \dot{\psi}_\theta^t [J_\lambda^{-1} + \lambda J_\lambda^{-1} [(I + \lambda \tau_0^2) \dot{\psi}_\theta \dot{\psi}_\theta^t] J_\lambda^{-1}] \dot{\psi}_\theta$ .

## Models with Local Misspecification - Exponential or Weibull?

Let  $Y_1, \dots, Y_n$  be i.i.d. **Weibull** with parameters  $\theta_0 = 0.34$  and  $\gamma = 1 + \delta/\sqrt{n}$ .

Note that  $\delta = 0$  is the **exponential distribution** (narrow).

Comparing the exponential (narrow), Weibull (wide) and **FRL** at estimating

$$\psi(\theta, \gamma) = \Pr\{Y_1 > 1\}.$$

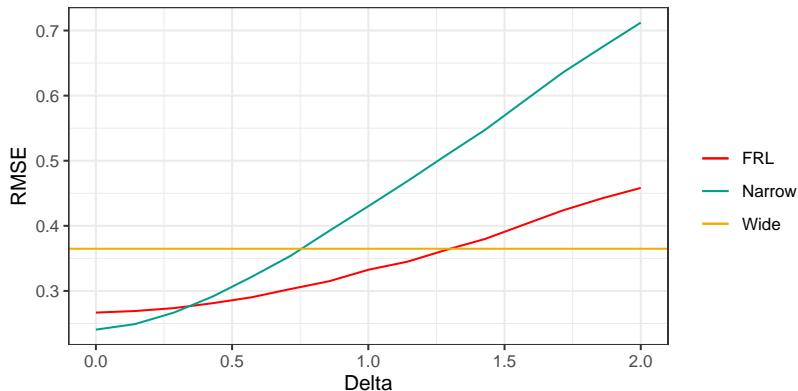
## Models with Local Misspecification - Exponential or Weibull?

Let  $Y_1, \dots, Y_n$  be i.i.d. **Weibull** with parameters  $\theta_0 = 0.34$  and  $\gamma = 1 + \delta/\sqrt{n}$ .

Note that  $\delta = 0$  is the **exponential distribution** (narrow).

Comparing the exponential (narrow), Weibull (wide) and **FRL** at estimating

$$\psi(\theta, \gamma) = \Pr\{Y_1 > 1\}.$$



## Focused regularised regression

The general idea and framework is easily extended to regression models.

## Focused regularised regression

The general idea and framework is easily extended to regression models.

A canonical FRL construction is then

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{ \hat{\mu}(z_i) - z_i^t \beta \}^2 \right\},$$

where  $\hat{\mu}$  is an **alternative estimate of the mean** and  $I$  is a set of important and/or control ‘individuals’.

## Focused regularised regression

The general idea and framework is easily extended to regression models.

A canonical FRL construction is then

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{\hat{\mu}(z_i) - z_i^t \beta\}^2 \right\},$$

where  $\hat{\mu}$  is an **alternative estimate of the mean** and  $I$  is a set of important and/or control ‘individuals’.

**Example:** A simple model combined with a **sophisticated non-parametric model** for  $\hat{\mu}(\cdot)$ ; e.g. to control for missing interactions.

## Focused regularised regression

The general idea and framework is easily extended to regression models.

A canonical FRL construction is then

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{ \hat{\mu}(z_i) - z_i^t \beta \}^2 \right\},$$

where  $\hat{\mu}$  is an **alternative estimate of the mean** and  $I$  is a set of important and/or control ‘individuals’.

**Example:** A simple model combined with a **sophisticated non-parametric model** for  $\hat{\mu}(\cdot)$ ; e.g. to control for missing interactions.

**Example:** To integrate an estimated model with the output from a **physical or mechanistic** model; weather, hydrology, biology, . . .

## Focused regularised regression

The general idea and framework is easily extended to regression models.

A canonical FRL construction is then

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{ \hat{\mu}(z_i) - z_i^t \beta \}^2 \right\},$$

where  $\hat{\mu}$  is an **alternative estimate of the mean** and  $I$  is a set of important and/or control ‘individuals’.

**Example:** A simple model combined with a **sophisticated non-parametric model** for  $\hat{\mu}(\cdot)$ ; e.g. to control for missing interactions.

**Example:** To integrate an estimated model with the output from a **physical or mechanistic** model; weather, hydrology, biology, . . . .

**Example:** To integrate local data with external data where we **do not have access to raw data**, relying on an estimated  $\hat{\mu}(\cdot)$  for integration.



## Focused regularised regression

The general idea and framework is easily extended to regression models.

A canonical FRL construction is then

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{ \hat{\mu}(z_i) - z_i^t \beta \}^2 \right\},$$

where  $\hat{\mu}$  is an **alternative estimate of the mean** and  $I$  is a set of important and/or control ‘individuals’.

**Example:** A simple model combined with a **sophisticated non-parametric model** for  $\hat{\mu}(\cdot)$ ; e.g. to control for missing interactions.

**Example:** To integrate an estimated model with the output from a **physical or mechanistic** model; weather, hydrology, biology, . . . .

**Example:** To integrate local data with external data where we **do not have access to raw data**, relying on an estimated  $\hat{\mu}(\cdot)$  for integration.

**Example:** To regularise a complex model, by penalising towards a **simple model** where data is sparse.

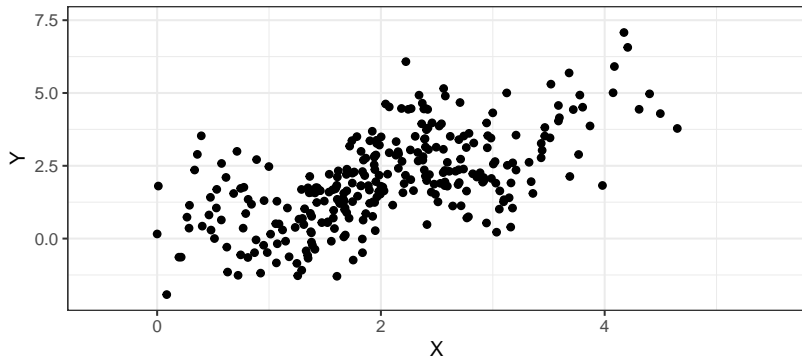
## Illustration - Focused regularised regression

Simulated data with

$$Y_i = \mu(x_i) + \epsilon_i$$

for **some smooth function**  $\mu(\cdot)$  and independent  $\epsilon_1, \dots, \epsilon_n$ .

Simulated Data



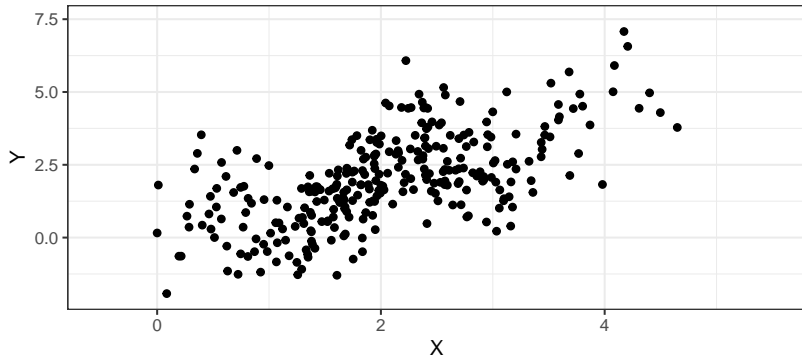
## Illustration - Focused regularised regression

Simulated data with

$$Y_i = \mu(x_i) + \epsilon_i$$

for **some smooth function**  $\mu(\cdot)$  and independent  $\epsilon_1, \dots, \epsilon_n$ .

Simulated Data



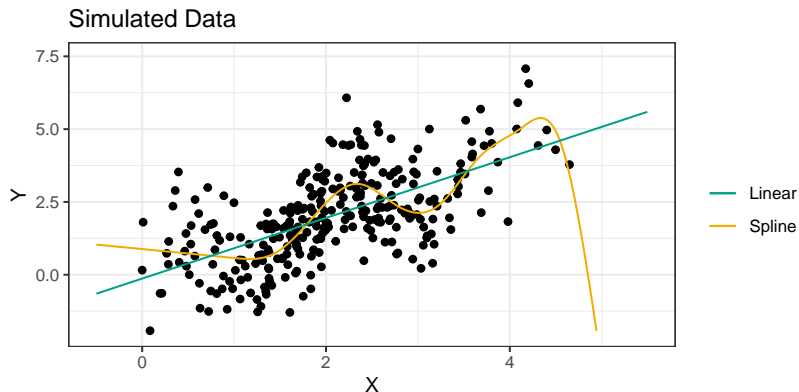
A model that captures both the **overall trend** and **the detailed behaviour**?

Useful for extrapolation.

## Illustration - Focused regularised regression

A **smooth spline** effectively capture detailed behaviour where data is dense

Can use a **simple linear model** to capture the overall trend.



How to combine?

## Illustration - Focused regularised regression

Inspired by the FRL setup

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(s_\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{(\hat{a} + \hat{b}z_i) - s_\beta(z_i)\}^2 \right\},$$

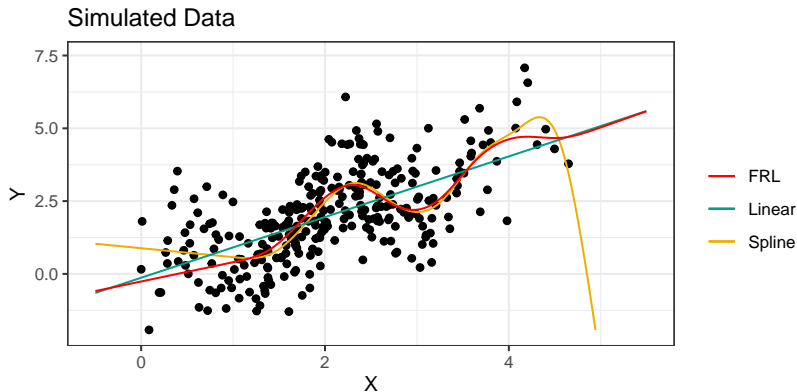
where  $s_\beta$  is a smooth spline and  $I$  is a set of control points – some below  $x = 0$  and some above  $x = 4$ .

## Illustration - Focused regularised regression

Inspired by the FRL setup

$$\hat{\theta}_\lambda = (\hat{\beta}_\lambda, \hat{\sigma}_\lambda) = \arg \max_{\beta, \sigma} \left\{ \ell_n(s_\beta, \sigma) - \frac{1}{2} \lambda n \frac{1}{|I|} \sum_{z_i \in I} \{(\hat{a} + \hat{b}z_i) - s_\beta(z_i)\}^2 \right\},$$

where  $s_\beta$  is a smooth spline and  $I$  is a set of control points – some below  $x = 0$  and some above  $x = 4$ .



## Concluding Remarks

Just do it.

A straightforward method for **improving robustness** of parametric models.

And can make **inference more focused**.

**Large-sample theory** justify the use in simple models.

Works well for **regression** models.

And stationary **time series**.

**Bootstrapping** techniques also works well (in simulated data examples).

Link to **empirical likelihood** and **empirical Bayes**.