

# Efficient statistical inference for high dimensional and nonparametric models

Natalia Bochkina

School of Mathematics and Maxwell Institute

N.Bochkina@ed.ac.uk

13 October 2014

# High dimensional data and modern statistics

Availability of **noisy high dimensional data** from biology and medicine (genomics, genetics, tomography, brain imaging), ecology, social networks and other data (netflix problem) triggered development of “**high  $p$  small  $n$** ” paradigm in statistics where the number of unknowns  $p$  is higher than the sample size  $n$ .

**Classical statistical inference** works if  $p$  is fixed as  $n \rightarrow \infty$ .

Current data has  $p \rightarrow \infty$  as  $n \rightarrow \infty$ , often  $p/n \rightarrow \text{const}$  or  $p/n \rightarrow \infty$ .

To ensure **consistent and efficient inference**, this required development of novel statistical methods, often embedding **a priori information** elicited from experts.

- **Main statistical methods**: penalised likelihood and Bayesian models
- **Computational challenges**: to implement the methods efficiently for a large number of unknowns
- **Mathematical challenges**: to come up with statistical inference methods that guarantee consistency and efficiency

It is often of interest to recover structure in these types of data.

# High dimensional data and modern statistics

Availability of **noisy high dimensional data** from biology and medicine (genomics, genetics, tomography, brain imaging), ecology, social networks and other data (netflix problem) triggered development of “**high  $p$  small  $n$** ” paradigm in statistics where the number of unknowns  $p$  is higher than the sample size  $n$ .

**Classical statistical inference** works if  $p$  is fixed as  $n \rightarrow \infty$ .

Current data has  $p \rightarrow \infty$  as  $n \rightarrow \infty$ , often  $p/n \rightarrow \text{const}$  or  $p/n \rightarrow \infty$ .

To ensure **consistent and efficient inference**, this required development of novel statistical methods, often embedding **a priori information** elicited from experts.

- **Main statistical methods:** penalised likelihood and Bayesian models
- **Computational challenges:** to implement the methods efficiently for a large number of unknowns
- **Mathematical challenges:** to come up with statistical inference methods that guarantee consistency and efficiency

It is often of interest to recover structure in these types of data.

# High dimensional data and modern statistics

Availability of **noisy high dimensional data** from biology and medicine (genomics, genetics, tomography, brain imaging), ecology, social networks and other data (netflix problem) triggered development of “**high  $p$  small  $n$** ” paradigm in statistics where the number of unknowns  $p$  is higher than the sample size  $n$ .

**Classical statistical inference** works if  $p$  is fixed as  $n \rightarrow \infty$ .

Current data has  $p \rightarrow \infty$  as  $n \rightarrow \infty$ , often  $p/n \rightarrow \text{const}$  or  $p/n \rightarrow \infty$ .

To ensure **consistent and efficient inference**, this required development of novel statistical methods, often embedding **a priori information** elicited from experts.

- **Main statistical methods**: penalised likelihood and Bayesian models
- **Computational challenges**: to implement the methods efficiently for a large number of unknowns
- **Mathematical challenges**: to come up with statistical inference methods that guarantee consistency and efficiency

It is often of interest to recover structure in these types of data.

# Statistical inference for high dimensional data

**Likelihood:**  $\hat{Y} = (Y_1, \dots, Y_n) \sim f(Y | \theta)$ , for some  $\theta \in \Theta \subseteq \mathbb{R}^p$ ,  $p \gg n$ .

Aim: to estimate unknown  $\theta$ , its confidence region, make decisions.

**Penalised log likelihood estimator:**

$$\hat{\theta} = \arg \min_{\hat{\theta}} \left[ -\log p(Y | \hat{\theta}) + \text{pen}(\hat{\theta}) \right]$$

where the penalty reflects desirable properties of the solution, e.g. sparsity.

**Problems:**

- Construction of confidence regions for  $\hat{\theta}$  and other decision making.
- Assumptions are often not verifiable.

**Bayesian model:**

Given prior  $p(\theta)$ , **posterior distribution** is

$$p(\theta | Y) = \frac{f(Y | \theta) p(\theta)}{\int_{\Theta} f(Y | \theta) p(\theta) d\theta},$$
$$\hat{\theta} = \arg \max_{\hat{\theta}} \mathbb{E} \left( d(\hat{\theta}, \theta) | Y \right)$$

given a distance  $d$  on  $\Theta$ . Bayesian analogues of a confidence region and decision making can be constructed.

# Challenges

- **Computational**: construct a fast algorithm to compute  $\hat{\theta}$  (and  $p(\theta | Y)$ ).
- **Mathematical**: choose a penalty  $pen(\theta)$  / prior  $p(\theta)$  to ensure that  $\hat{\theta}$  is **consistent and efficient**, i.e.

$$\mathbb{E}[d(\hat{\theta}, \theta)]^2 \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

at the best possible rate.

**For Bayesian inference**, efficiency is related to local concentration of the posterior distribution around the true value of  $\theta$  (Bernstein-von Mises theorem).

# Research problems

- **Modelling and decision making for genomic data**, including model checks  
Bochkina et al (2006, 2007, 2010)
- **Modelling dependence structure in genomics data and data integration**  
A. Caballe, Bochkina, C.-D. Meyer
- **Statistical inference for compound sparse high dimensional problems**  
Bochkina & Ritov (2011)
- **Concentration of posterior distribution (Bernstein – von Mises theorem) for nonregular and misspecified models, with application to tomography**  
Bochkina and Green (2014), Bochkina (2013)
- **Concentration of posterior distribution for semiparametric models with functional nuisance parameter**  
Bochkina and Rousseau
- **Adjusting Bayesian inference for approximated models**

## Concentration of posterior distribution (Bernstein–von Mises theorem)

For **correctly specified regular models**, as  $n \rightarrow \infty$ ,

$$l_{\theta_{\text{true}}}^{1/2}(\theta - \theta_{\text{true}}) \mid \mathbf{Y} \sim N_p(\Delta_n, l_p)$$

where  $l_{\theta_{\text{true}}}$  is Fisher information:

$$l_{\theta_{\text{true}}} = \mathbb{E}_{\theta_{\text{true}}} \left( \frac{\partial \log f(\mathbf{Y} \mid \theta_{\text{true}})}{\partial \theta} \right)^2.$$

That is, Bayesian inference in asymptotically optimal in frequentist sense.

For iid models,  $l_{\theta_{\text{true}}} = ni_{\theta_{\text{true}}}$  leading to standard  $\sqrt{n}$  parametric convergence rate.

For **nonregular models**, where  $\theta_{\text{true}}$  or  $\theta^*$  are on the boundary of  $\Theta$ ,

$$n(\theta - \theta^*)_j \mid \mathbf{Y} \sim \Gamma(\alpha, b_j) \quad \text{as } n \rightarrow \infty$$

for some directions  $j$  where  $\alpha$  is parameter of the prior distribution.

For **misspecified regular models**, where  $f_{\text{true}}(\mathbf{Y}) \notin \{f(\mathbf{Y} \mid \theta), \theta \in \Theta\}$ ,

$$V_{\text{true}}^{1/2}(\theta - \theta_{\text{true}}) \mid \mathbf{Y} \sim N_p(\Delta_n, l_p)$$

where

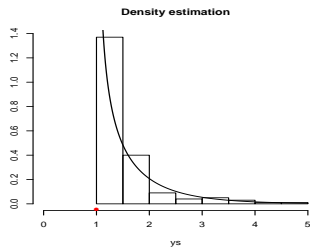
$$V_{\text{true}} = \mathbb{E}_{\text{true}} \left( \frac{\partial \log f(\mathbf{Y} \mid \theta^*)}{\partial \theta} \right)^2,$$

and  $\theta^*$  corresponds to the model  $f(\mathbf{Y} \mid \theta^*)$  closest to the true one  $f_{\text{true}}(\mathbf{Y})$  in Kullback-Leibler distance.



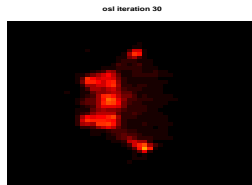
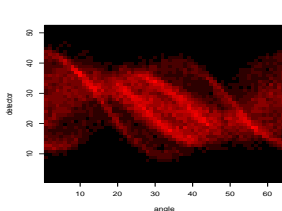
# Application to nonparametric and inverse problems, tomography

- **Density estimation:**  $(Y_1, \dots, Y_n)$  is a sample from density  $f$ .
- **Nonparametric regression:**  $(Y_1, \dots, Y_n)$  are noisy observations of  $f$  at points  $(t_1, \dots, t_n)$ .
- **Inverse problem:**  $(Y_1, \dots, Y_n)$  are noisy observations of  $\mathcal{A}(f)$  (indirect observations of  $f$ ) at points  $(t_1, \dots, t_n)$ .



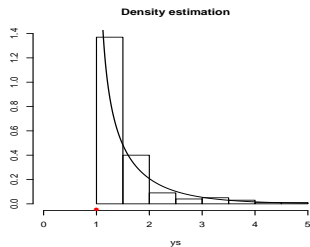
**Aim:** estimate unknown function  $f$ .

Inverse problem, tomography (plots by P.J.Green)



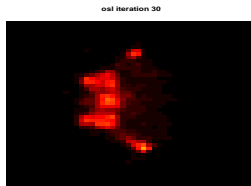
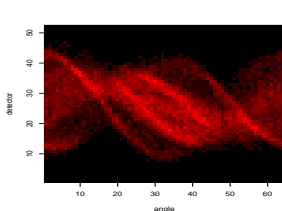
# Application to nonparametric and inverse problems, tomography

- **Density estimation:**  $(Y_1, \dots, Y_n)$  is a sample from density  $f$ .
- **Nonparametric regression:**  $(Y_1, \dots, Y_n)$  are noisy observations of  $f$  at points  $(t_1, \dots, t_n)$ .
- **Inverse problem:**  $(Y_1, \dots, Y_n)$  are noisy observations of  $\mathcal{A}(f)$  (indirect observations of  $f$ ) at points  $(t_1, \dots, t_n)$ .



**Aim:** estimate unknown function  $f$ .

## Inverse problem, tomography (plots by P.J.Green)



## Adjusting Bayesian inference for approximated models

For complex Bayesian models, **approximate models** are often fitted to speed up the computation.

They often **underestimate uncertainty** in the posterior distribution.

The idea is

- to view approximate models as misspecified models
- to use BvM for misspecified models to adjust the inference

## Adjusting Bayesian inference for approximated models

For complex Bayesian models, **approximate models** are often fitted to speed up the computation.

They often **underestimate uncertainty** in the posterior distribution.

The idea is

- to view approximate models as misspecified models
- to use BvM for misspecified models to adjust the inference