



# Project C4

## Regression approaches for large-scale high-dimensional data

Prof. Dr. Katja Ickstadt, Prof. Dr. Christian Sohler

Problem

- ▶ Problem: Large number of observations or large dimension
- ▶ Consequence: Modern regression approaches reach limits of scalability
- Goal: Develop highly efficient regression approaches for massive data**

### Bayesian Regression

$$\pi(\beta | X, Y) \propto \mathcal{L}(Y | X, \beta) \cdot \pi(\beta)$$

- ▶  $\ell_p$  regression
- ▶ Hierarchical prior models

### Generalised Linear Models (GLMs)

$$h(\mathbb{E}(Y)) = X\beta.$$

Examples:

- ▶ Poisson regression for count data
- ▶ Logistic regression for binary data

### Structural Constraints

$$\min_{\xi \in \mathcal{N}} \|\xi - Y\|, \text{ for some } \mathcal{N} \subseteq \mathbb{R}^n$$

- ▶ E.g. monotonicity, unimodality
- $\mathcal{N} = \{\xi \mid \exists i \in [n]: \xi_1 \leq \dots \leq \xi_i \geq \dots \geq \xi_n\}$
- ▶ Application: peak detection

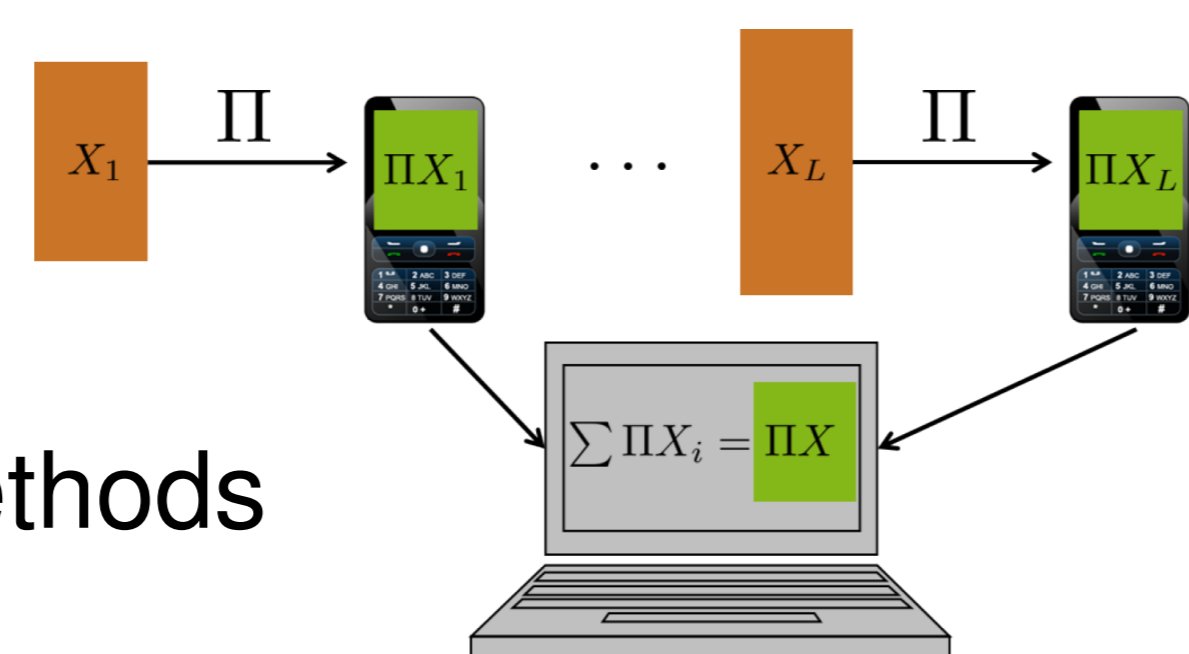
### Dimensionality reduction

- ▶ Feature selection
- ▶ Variable interactions
- ▶ SNP and genome-wide data

Methodology

### Data Reduction with Guaranteed Little Distortion

- ▶ Streaming
- ▶ Distributed
- ▶ Succinct representation



### Sampling and projection methods

- ▶  $\epsilon$ -subspace embeddings
- ▶  $\epsilon$ -coresets
- ▶  $\Pi \in \mathbb{R}^{k \times n}, k \ll n$ :

$$\forall \beta \in \mathbb{R}^d: (1 - \epsilon)\|X\beta\| \leq \|\Pi X\beta\| \leq (1 + \epsilon)\|X\beta\|$$

### Algorithmic Principle

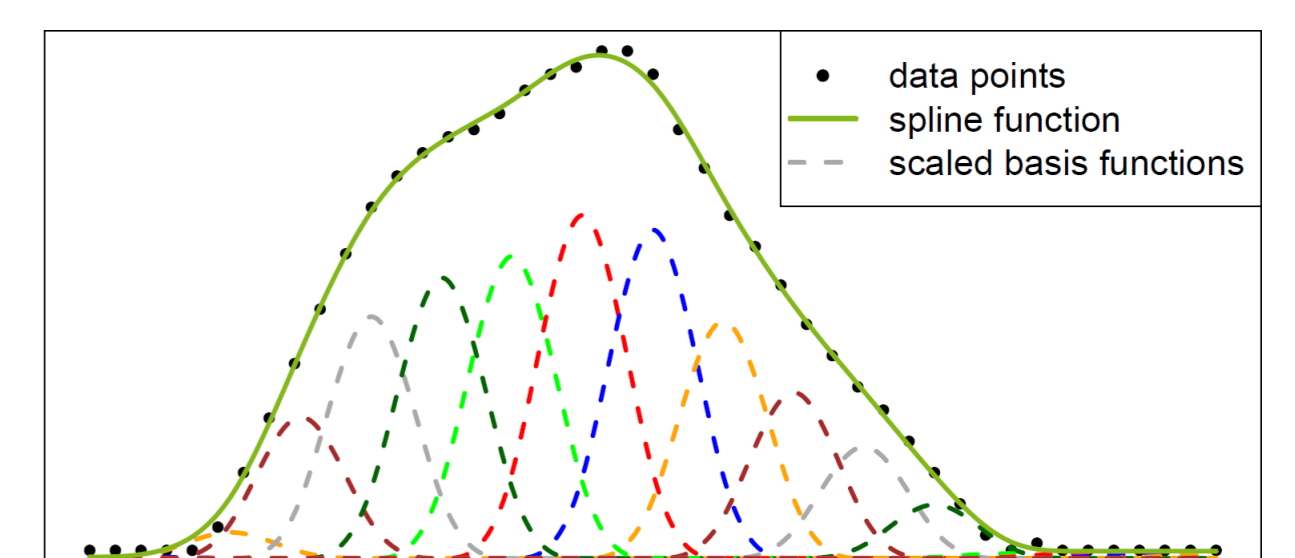
1. Data reduction
2. Statistical analysis

$$[X, Y] \xrightarrow{\Pi} [\Pi X, \Pi Y]$$

$$\pi(\beta | X, Y) \approx_{\epsilon} \pi(\beta | \Pi X, \Pi Y)$$

### Spline Regression Model Structural Constraints

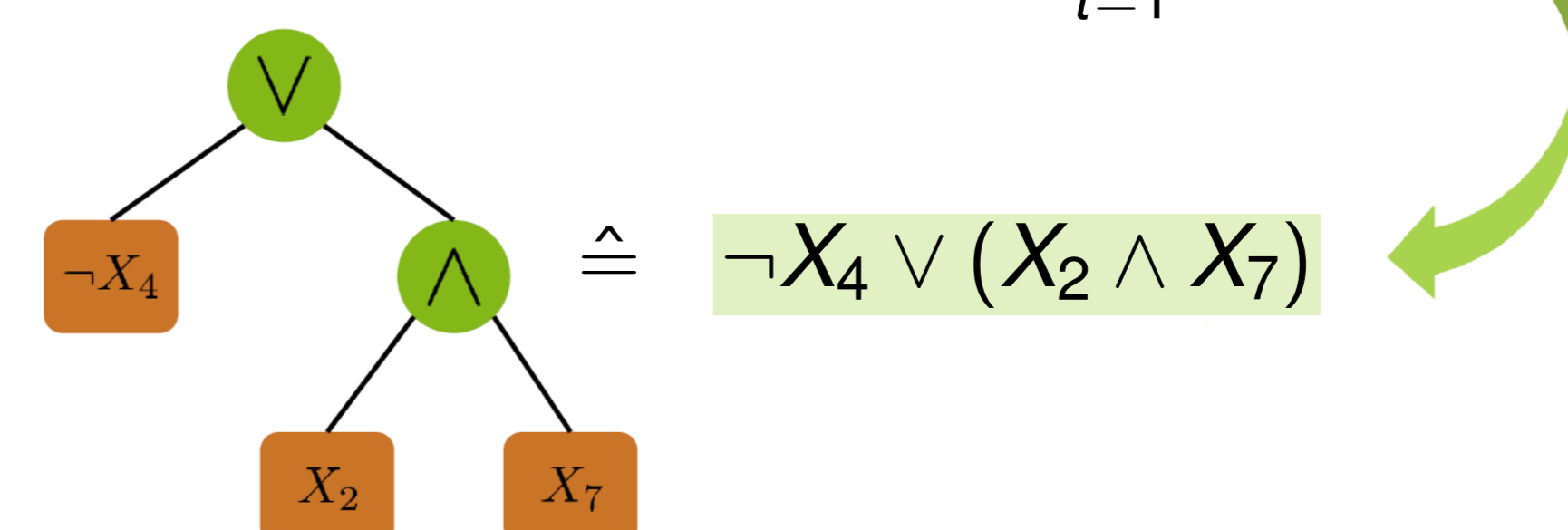
$$s(x) = \sum_{i=1}^B \beta_i b_i(x)$$



### Logic regression

- ▶ Interactions of (binary) variables
- ▶ Variable importance measures

$$h(\mathbb{E}(Y)) = \beta_0 + \sum_{t=1}^T \beta_t L_t$$



Results

### Scalable Bayesian Regression

- ▶  $\epsilon$ -subspace embeddings preserve
  - ▶  $\ell_p$ -generalised linear regression,  $p \in [1, \infty)$
  - ▶ Hierarchical prior models

### GLMs

- ▶ Lower bounds  $k \in \Omega(n/\log n)$

### Poisson regression

- ▶ Statistical model relaxation yields data summary
- ▶ Close approximation of maximum likelihood estimator

B4

### Logistic Regression

- ▶ Geometric parameter captures difficulty of reduction

$$\mu(X) = \sup_{\beta \in \mathbb{R}^d \setminus \{0\}} \frac{\|(X\beta)^-\|_1}{\|(X\beta)^+\|_1}$$

A2

- ▶  $\mu$ -complex data admits coreset of size  $O((\mu d \log n/\epsilon)^{O(1)})$

[Geppert, Ickstadt, Munteanu, ..., Sohler, Stat Comput 2017]  
 [Müller, MT, 2016]  
 [Munteanu, Diss, 2018]  
 [Molina, Munteanu, Kersting, AAAI, 2018]  
 [Munteanu, Schwiigelshohn, Sohler, ..., NIPS, 2018]

### Spline Regression

- ▶ Able to model unimodality
- ▶ Limited number  $B$  of splines (succinct representation)
- ▶ Peak detection in IMS spectra

B1

### Scalable Logic Regression

- ▶ Useful importance measures *cross-leverage-scores*:

$$cl_i = U_i^T U_{d+1}, \text{ for } i \in [d],$$

where  $U$  is an orthonormal basis for  $[X, Y]$ .

- ▶ Subsampling retains interesting variables
- ▶ Main effects and interactions
- ▶ Logic regression operates efficiently on reduced variable set

[Köllmann, Diss, 2016]  
 [Wollenberg, MT, 2016]