

Predicting from high-dimensional molecular data and environmental variables in stratified samples

Norbert Krautenbacher, Christiane Fuchs, Fabian J. Theis
Institute of Computational Biology, Helmholtz Zentrum München, Department of Mathematics, Technische Universität München

Motivation: Prediction of childhood asthma

Environment + Genetics → Asthma / No asthma

→ **Objective:** Build risk score for childhood asthma based on genetic and environmental predictors

→ **Challenges:**

- $n \ll p$
- biased sample
- combining high-dim. with low-dim. data

→ Can childhood asthma be predicted by genetics and environment?

Correcting for the sample bias

Survey design (simplified)

- Interest in **influence of farm exposure on asthma**
- For cost reasons: detailed investigations only on **subsample of finite population**

→ How to take sample?

→ **Enrichment of informative observations by Stratified two-phase random selection process** (simplified, let $N=10000$):

Population of size $N=10000$

Group	Stratum	Count
Farm Children	Asthma	300
	Controls	3000
Nonfarm Children	Asthma	1800
	Controls	4900

1st selection → 2nd selection → sample of size $n=400$ (= given dataset)

→ **Given dataset is result of selection process**

→ **Correction for sample bias necessary**

→ otherwise biased estimates

Data

- Data from GABRIELA study
- Project in collaboration with Dr von Mutius, Hauersches Kinderspital

	~40 environmental predictors			2.5 million genetic predictors			Outcome
	Farm	Straw	Stable	SNP1	SNP2	SNP3	Asthma
Child 1	0	0	0	0.45	1	0.21	1
Child 2	1	1	0	1.76	2	1.20	0
Child 3	1	0	1	1.01	0	0.94	1
...
Child 1707	0	1	1	0.22	1	1.99	0

→ $n \ll p$ problem

→ **dimension reduction necessary**

Results

• **SNPs only:** weak prediction power

• **Environment only:** better prediction than for genetics

• **Combination:** no improvement

Performance of several methods for asthma on SNPs from literature research + environment - weighted validation

Method: Log-Regr, Log-Regr+weights, LASSO, LASSO+weights, PCA+Log-Regr, PCA+RF, PCA+SNPs+Log-Regr, PCA+SNPs+RF, RF

predictors: 19 SNPs (green), 328 SNPs (cyan)

Approaches

Correcting for Sample Bias:

- In **training procedure:** algorithm-specific **weighting approaches**
- In **testing procedure:** **Validation** by unbiased empirical estimate of the loss calculated by the **weighted mean over losses** per stratum

$$l_u = \frac{1}{m} \sum_{s=1}^m w_s \cdot l(h(x_s), y_s)$$

$n \ll p$ approaches:

- Reducing dimension of SNP data by
 - Manual selection of **SNPs by literature research**
 - **Univariate feature selection** by survey logistic regression
- Application of appropriate **Learning procedures:**
 - LASSO
 - Random Forest
 - PCA-based approaches

Conclusions

- **Best prediction** by PCA followed by RF
- **Better prediction on environment** compared to genetics
- **Combination** of predictors **doesn't improve** prediction → environmental effect covers genetic effect

Outlook

Improve prediction by

- **Combining genetic with environmental data** by special approaches
- Incorporation of sample bias
- Taking **all SNPs** into account in further **multivariate methods**

• Ege et al. Gene-environment interaction for childhood asthma and exposure to farming in Central Europe, *JACI*, 2011

• R. De Bin, W. Sauerbrei, A.-L. Boulesteix Investigating the prediction ability of survival models based on both clinical and omics data: two case studies, *Statistics in Medicine*, 2014