

Experimental design for genome-wide association studies

Christoph Lippert, Oliver Stegle, Hannes Nickisch, Karsten Borgwardt, Detlef Weigel

Max Planck Institutes, Tübingen

Approximate Bayesian learning

We follow the ideas developed in [Nickisch and Seeger, 2009] and approximate the intractable exact posterior, $P(\theta|\mathbf{y}, \mathbf{X}; \sigma^2, \tau)$, by a **Gaussian approximation** of the form

 $\mathcal{Q}_{\gamma}(\boldsymbol{\theta}; \sigma^2, \tau) = \mathcal{N}(\mathbf{y} | \mathbf{X} \boldsymbol{\theta}, \sigma^2 \mathbf{I}) \prod t_i(\theta_i, \gamma_i),$

with variational parameters $\boldsymbol{\gamma} \in \mathbb{R}^d_+$. $t_i(\theta_i, \gamma_i) \propto \mathcal{N}(0, \boldsymbol{\gamma}_i)$ are Gaussian site functions that **lower bound** the exact Laplace sites. The product of the likelihood term and the approximate sites in Equation 4 is tractable and can be written as

$$\mathcal{Q}_{\gamma}(\boldsymbol{\theta}; \sigma^2, \tau) = \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}_{\boldsymbol{\theta}},$$

with mean μ_{θ} and covariance matrix Σ_{θ} ,

$$\boldsymbol{\mu}_{\boldsymbol{\theta}} = \mathbf{A}^{-1} \mathbf{X}^{\mathrm{T}} \mathbf{y}, \qquad \boldsymbol{\Sigma}_{\boldsymbol{\theta}} = \sigma^2 \mathbf{A}^{-1}, \qquad \mathbf{A}$$

Under this approximation, the prediction at an unseen test input \mathbf{x}_{\star} again has a Gaussian distribution,

$$\mathbf{y}_{\star} \sim \mathcal{N}(\boldsymbol{\mu}_{\star}, \sigma_{\star}^2), \qquad \boldsymbol{\mu}_{\star} = \mathbf{x}_{\star} \boldsymbol{\mu}_{\boldsymbol{\theta}},$$

The variational parameters γ are determined such that \mathcal{Q}_{γ} approximates the exact posterior well. We developed an algorithm that minimises the convex relaxation of the KL-divergence between the exact and approximate posterior from [Nickisch and Seeger, 2009] in runtime that is **linear** in the number of SNPs, $O(dn^3)$.

Experimental design

We perform **greedy blockwise experimental design** using alternative selection criteria.

Mean marginal entropy

We use the expected reduction in mean marginal entropy of the posterior as a design criterion.

$$\Delta H_{\rm mm} = H_{\rm mm} - H_{\rm mm}^{\mathcal{C}}.$$
 (8)

The reduction in mean marginal entropy ΔH_{mm} of a fully trained model, defined by the approximate posterior of $\boldsymbol{\theta}$, is

$$\frac{1}{2} \sum_{d=1}^{D} \log \boldsymbol{\Sigma}_{\boldsymbol{\theta}}[d, d] - \frac{1}{2} \sum_{d=1}^{D} \log \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{\mathcal{C}}[d, d].$$
(9)

Population posterior variance

Another experimental design criterion is the expected posterior variance evaluated on a target population \mathcal{P}_{\cdot} $\sum_{\mathbf{x}_{\mathcal{P}} \in \mathcal{D}} \mathbf{x}_{\mathcal{P}} \Sigma_{\boldsymbol{\theta}}^{\mathcal{C}} \mathbf{x}_{\mathcal{P}}^{\mathrm{T}}.$

.

position

Note that \mathcal{P} usually is either the set of *all genotypes* or in a transductive setting an *independent test*

In both criteria the **posterior covariance matrix** $\Sigma_{\boldsymbol{\rho}}^{\mathcal{C}}$ including a number of candidate instances $\mathbf{X}_{\mathcal{C}}$ is approximated by

$$\boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{\mathcal{C}} = \sigma^2 \left(\mathbf{X}_{\mathcal{C}}^{\mathrm{T}} \mathbf{X}_{\mathcal{C}} + \mathbf{A} \right)$$



(11)

- 166 accessions



Discussion

- per iteration.

Outlook

- Experimental design on 1200 recently genotyped A. thaliana accessions to extend the study by [Atwell et al., 2010].
- Optimal solution of blockwise selection
- How to determine the *optimal block size*?
- Extend to **multiple phenotypes**
- Correct for **population structure** Non-i.i.d. data
 - Remove spurious associations

References

- http://www.1000genomes.org.
- inbred lines. Nature, 2010.
- 851-861, 2007.

• By doing **greedy blockwise** experimental design it is feasable to include *multiple individuals*

Individuals

1000 Genomes Project. A deep catalog of human genetic variation. World Wide Web electronic publication, 2009. URL

1001 Genomes. A catalog of *Arabidopsis thaliana* genetic variation, 2009. URL http://www.1001genomes.org.

S. Atwell, Y. Huang, B. Vilhálmsson, G. Willems, M. Horton, Y. Li, D. Meng, A. Platt, A. Tarone, T. Hu, R. Jiang, N. Muliyati, X. Zhang, M. Amer, I. Baxter, B. Brachi, J. Chory, C. Dean, M. Debieu, J. de Meaux, J. Ecker, N. Faure, J. Kniskern, J. Jones, T. Michael, A. Nemri, F. Roux, D. Salt, C. Tang, M. Todesco, M. Traw, D. Weigel, P. Marjoram, J. Borevitz, J. Bergelson, and M. Nordborg. Genome-wide association study of 107 phenotypes in Arabidopsis thaliana

H. Nickisch and M. Seeger. Convex variational bayesian inference for large scale generalized linear models. In *Proceedings* of the 26th Annual International Conference on Machine Learning, pages 761–768, 2009. The International HapMap Consortium. A second generation human haplotype map of over 3.1 million snps. *Nature*, 449: