

Causal stability ranking

Daniel J. Stekhoven^{1,2,3,*}, Izabel Moraes⁴, Gardar Sveinbjörnsson¹, Lars Hennig^{4,5}, Marloes H. Maathuis¹ and Peter Bühlmann^{1,3}

¹Seminar for Statistics, Department of Mathematics, ²Life Science Zurich PhD Program on Systems Biology of Complex Diseases, ETH Zurich, 8092 Zurich, Switzerland, ³Competence Center for Systems Physiology and Metabolic Diseases, 8092 Zurich, Switzerland, ⁴Uppsala BioCenter, Department of Plant Biology and Forest Genetics, Swedish University of Agricultural Sciences and Linnean Center for Plant Biology, 750 07 Uppsala, Sweden and ⁵Plant Biotechnology, Department of Biology, ETH Zurich, 8092 Zurich, Switzerland

Associate Editor: Olga Troyanskaya

ABSTRACT

Genotypic causes of a phenotypic trait are typically determined via randomized controlled intervention experiments. Such experiments are often prohibitive with respect to durations and costs, and informative prioritization of experiments is desirable. We therefore consider predicting stable rankings of genes (covariates), according to their total causal effects on a phenotype (response), from observational data. Since causal effects are generally non-identifiable from observational data only, we use a method that can infer lower bounds for the total causal effect under some assumptions. We validated our method, which we call Causal Stability Ranking (CStaR), in two situations. First, we performed knock-out experiments with *Arabidopsis thaliana* according to a predicted ranking based on observational gene expression data, using flowering time as phenotype of interest. Besides several known regulators of flowering time, we found almost half of the tested top ranking mutants to have a significantly changed flowering time. Second, we compared CStaR to established regression-based methods on a gene expression dataset of *Saccharomyces cerevisiae*. We found that CStaR outperforms these established methods. Our method allows for efficient design and prioritization of future intervention experiments, and due to its generality it can be used for a broad spectrum of applications.

Availability: The full table of ranked genes, all raw data and an example R script for CStaR are available from the Bioinformatics website.

Contact: stekhoven@stat.math.ethz.ch

Supplementary Information: Supplementary data are available at Bioinformatics online.

Received on March 22, 2012; revised on July 5, 2012; accepted on August 17, 2012

1 INTRODUCTION

The growing interest in causal inference (e.g. Kruglyak and Storey, 2009) has increased the need not only for methods able to handle this task but also for designed experimental validation. It is of general interest to infer the genotypic causes of a complex phenotypic trait (Glazier *et al.*, 2002). The classical approach relies on randomized controlled intervention experiments, e.g. knocking out a gene and observing the effect on the phenotype

relative to the wild-type organism. However, such intervention experiments are time consuming and expensive, and a prioritization with respect to most informative new experiments is very desirable. A genetic method to identify loci causing phenotypes or gene expression patterns is based on quantitative trait loci (QTL) and expression QTL (Gilad *et al.*, 2008; Kliebenstein, 2009). This can be a very powerful approach but it is limited to loci where genetic variation exists and to situations where segregating progeny of control crosses is available. Often, however, it is desirable to predict causal effects from purely observational data. We therefore consider the problem of predicting total causal effects from data obtained by observing a system without subjecting it to targeted interventions (observational data). This problem is generally ill-posed, but the recently proposed IDA method (Maathuis *et al.*, 2009, 2010) provides estimated lower bounds of total causal effects from observational data under some assumptions (Supplementary Section S1). However, these bounds come without a measure of uncertainty. We address this issue by introducing a new method combining IDA and a version of stability selection (Meinshausen and Bühlmann, 2010), which we call Causal Stability Ranking (CStaR; Fig. 1). The addition of stability selection to IDA provides two advantages. First, CStaR leads to a stable ranking of genes (covariates) according to the sizes of lower bounds for their predicted total causal effects, irrespective of the choice of the tuning parameter in stability selection. Second, under some additional assumptions, CStaR allows controlling an error rate of false-positive findings, namely the expected number of false positives and hence also the per-comparison error rate (PCER). CStaR results were confirmed in two biological scenarios using the simple model *Saccharomyces cerevisiae* and the more complex model *Arabidopsis thaliana*. Together, the built-in error measure and the success in finding relevant regulator genes make CStaR an excellent ranking method for the targeted design of experiments based on easily available resources.

2 METHODS

Based on observational training data and a set of required assumptions, CStaR predicts a lower bound for the total causal effect of a covariate on a response of interest, including a PCER for the false-positive selections. This is achieved by combining IDA (Section 2.1) with a version of stability selection (Section 2.2) on a range of different parameters. Predicted

*To whom correspondence should be addressed.

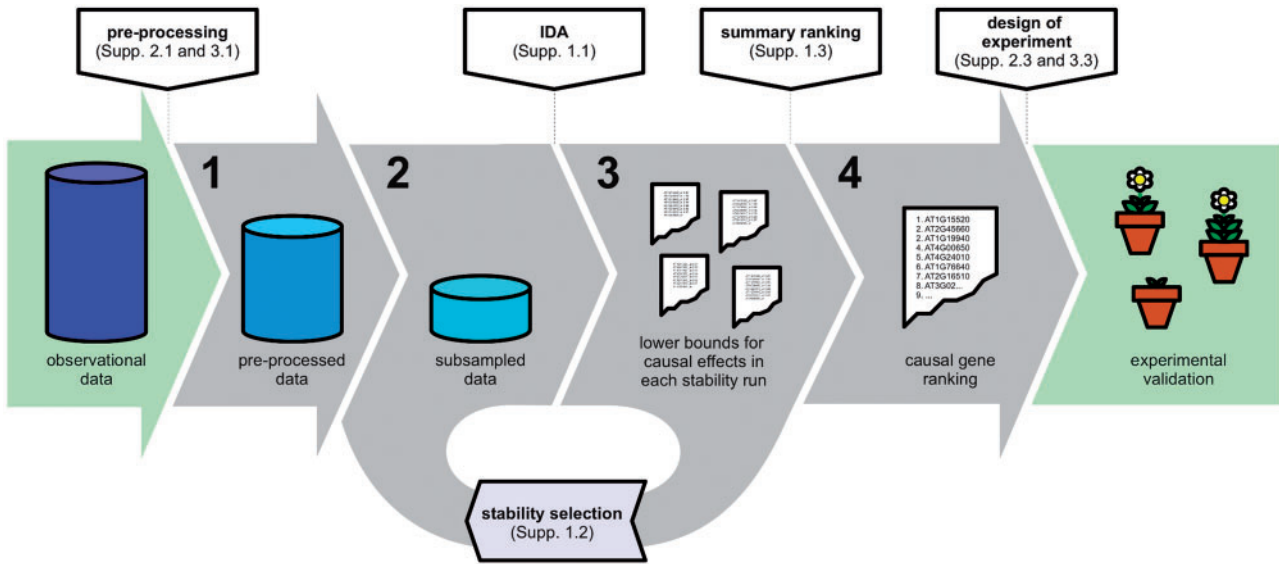


Fig. 1. Schematic overview of the methodological framework used in CStaR. After pre-processing the data (Step 1), lower bounds for the total causal effects are estimated 100 times using stability selection (Meinshausen and Bühlmann, 2010) according to the following procedure. A subsample of size $\lfloor n/2 \rfloor$ is repeatedly drawn from the total of n pre-processed data points (Step 2). On each subsample (or stability run), lower bounds for the total causal effects are estimated using IDA (Maathuis *et al.*, 2009) and used to rank the genes (Step 3, Section 2.1). Next, for a range of different q -values, we record the relative frequencies over the 100 stability runs that each gene appeared in the top q ranks (Section 2.2). The median rank over these different q s is used to generate the final ranking of the genes (Step 4). Furthermore, under additional assumptions, an upper bound for the PCER is estimated for each q -value and its corresponding relative frequency (Section 2.3). Finally, the gene ranking allows for design of new experiments. Thus, a biological validation using intervention experiments can be performed. We tested CStaR in two situations. First, on a publicly available compendium of 31 natural *A. thaliana* accessions consisting of $n = 47$ gene expression measurements, each with 21,326 genes and corresponding flowering time data (Lempe *et al.*, 2005; Supplementary Section S2.1). We performed biological intervention experiments according to the causal gene ranking (Table 1) by focusing on candidates that were not already known to control flowering time and for which mutant seeds were readily available (Supplementary Section S2.3). The biological experiments were analyzed using a two-sample Welch's t -test (Supplementary Section S2.4). The second validation was performed on a publicly available dataset in *S. cerevisiae* containing $n = 63$ observational and 234 interventional full-genome expression profiles, with $p = 5,361$ genes (Hughes *et al.*, 2000; Supplementary Section S3). Since this dataset includes both observational and interventional data, the validation was analyzed by comparing estimated total causal effects based on the observational data with inferred effects from the interventional data (Fig. 2)

total causal effects are ranked according to their stability aggregated over this range (Section 2.3).

2.1 Causal inference when the directed acyclic graph (DAG) is absent (IDA)

The IDA procedure (Maathuis *et al.*, 2009) is a statistical method that infers lower bounds for the absolute values of total causal effects on a response of interest from observational data under the assumption that the data come from an unknown DAG without hidden variables.

Suppose we have a dataset with n observations consisting of a response and p explanatory variables. Denoting by θ_j ($j = 1, \dots, p$), the true total causal effect of gene (covariate) j to the response (the total causal effect θ_j can be interpreted as follows: a change of gene j by one unit (one standard deviation) causes an average change of size θ_j in the response), the output of IDA is the estimated lower bound $\hat{\beta}_j$. It is shown (Maathuis *et al.*, 2009) that under certain assumptions (Supplementary Section 1) and as sample size n tends to infinity:

$$\hat{\beta}_j \xrightarrow{n \rightarrow \infty} \beta_j, \beta_j \leq |\theta_j|,$$

justifying the IDA procedure to infer lower bounds. These lower bounds are conservative: for example, if the lower bound is equal to zero, we would not make a statement that there is no causal effect (since the true total causal effect could be indeed equal to zero, or it could be larger than zero but the lower bound would not detect it). Based on the estimated lower bounds, we obtain a ranking of genes (covariates) with j_1 being the

index corresponding to the top rank, j_2 for the second best rank and so on:

$$\hat{\beta}_{j_1} \geq \hat{\beta}_{j_2} \geq \dots \geq \hat{\beta}_{j_p} \quad (1)$$

Under the assumption that the data come from an unknown DAG without hidden variables, the true total causal effect θ_j is generally non-identifiable from observational data, but lower bounds are. The conceptual idea for constructing lower bounds is as follows (Maathuis *et al.*, 2009). We first infer the so-called Markov equivalence class of all the DAGs (see Supplementary Section S1), which are compatible with the observational data. Using intervention calculus (Pearl, 2000), we derive all potential total causal effects based on each DAG G_r in the equivalence class (for every gene (covariate) j)

$$\{\theta_{j,r}; r = 1, \dots, m\} (j = 1, \dots, p),$$

and we define the true lower bounds as

$$\beta_j = \min_{r=1, \dots, m} |\theta_{j,r}| (j = 1, \dots, p). \quad (2)$$

Under our assumptions (see Supplementary Section S1), these (true) lower bounds β_j are identifiable from observational data, and the IDA algorithm yields the estimates $\hat{\beta}_j$ ($j = 1, \dots, p$). The main components of the IDA method are the PC-algorithm for estimating the Markov equivalence class of DAGs (Spirtes *et al.*, 2000) and a local algorithm for calculating the bounds β_j without enumerating all DAG members in the estimated Markov equivalence class (Maathuis *et al.*, 2009). It is

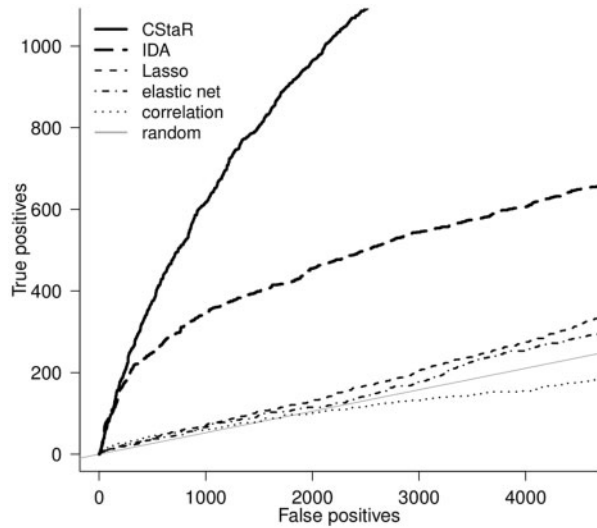


Fig. 2. True-positive selections (y-axis) versus false-positive selections (x-axis) for CStaR (solid) versus plain IDA (Maathuis *et al.*, 2009; long dashed), Lasso (Tibshirani, 1996; short dashed), elastic net (Zou and Hastie, 2005; dash dotted), the latter two using linear models and marginal correlation ranking (dotted) in the *S. cerevisiae* validation (Supplementary Section S3). Random guessing is indicated by the grey line. All methods were trained on the observational data. True positives were defined as the largest 5% of the effects (in absolute value) inferred from the interventional data

this local algorithm that makes the inference of these lower bounds based on thousands of genes (covariates) feasible. IDA is implemented in the R-package `pcalg` (Kalisch *et al.*, 2012).

2.2 Stability selection

CStaR incorporates a stability selection step (Meinshausen and Bühlmann, 2010). We draw 100 independent random subsamples of size $n/2$ and we run IDA on the subsampled data. In each subsampling run, which we also call stability run, we check whether gene (covariate) j has appeared among the top q variables when using the ranking as in equation (1) based on the subsampled data. We can then report the relative selection frequency $\hat{\Pi}_j$, among the 100 stability runs, that gene (covariate) j has appeared (or been selected) among the top q variables. These relative selection frequencies yield a stable list of genes (covariates): the index j_1 corresponds now to the most stably selected variable, and j_p to the least stable variable:

$$\hat{\Pi}_1 \geq \hat{\Pi}_2 \geq \dots \geq \hat{\Pi}_p. \quad (3)$$

Besides the increased stability in the ranking (3), stability selection is controlling the expected number of false-positive selections. Define the stably selected genes (covariates) as

$$\hat{S}_{\text{stable}} = \{j; \hat{\Pi}_j \geq \pi_{\text{thr}}\},$$

for some threshold $0.5 < \pi_{\text{thr}} \leq 1$. Denote the wrongly selected genes (false positives) by $V = |\hat{S}_{\text{stable}} \cap S_{\text{false}}|$, where S_{false} is the set of (false) genes (covariates) whose true lower bound $\beta_j = 0$, see (2). Then, for a given threshold π_{thr} and a given value of q [which influences (3)] we have, assuming an exchangeability condition (see Supplementary Section S1; Meinshausen and Bühlmann, 2010):

$$E[V] \leq \frac{1}{2\pi_{\text{thr}} - 1} \frac{q^2}{p} \quad (4)$$

and this leads to a bound for the PCER ($\text{PCER} = E[V]/p$). If a gene (covariate) j has relative selection frequency $\hat{\Pi}_j$, a bound for the corresponding PCER is given by

$$\frac{1}{2\hat{\Pi}_j - 1} \frac{q^2}{p^2}$$

2.3 Summary ranking

As novelty, we avoid choosing a specific q for the execution of stability selection by assessing the stability and the rank of each gene on a range of different q -values. This constitutes the main modification of the standard stability selection scheme and it also constitutes a useful simplification for the practitioner. This can be summarized graphically (Supplementary Fig. S1 gives an example for a single gene in the *A. thaliana* validation). We found that CStaR is relatively insensitive to the choice of the range of q s. However, down to a certain lower bound, small values of q lead to higher sensitivity and thus better results (see also Supplementary Section S3). If the q -values fall below such a lower bound, the ranking becomes unstable again. Finally, all genes are ranked according to the median rank with respect to the different q -values. Ties in the final ranking are sorted according to median total causal effect size.

2.4 Validation

We validated CStaR in two situations. First, we trained CStaR on a publicly available compendium of *A. thaliana* gene expression data and performed new biological validation experiments (Supplementary Section S2). The compendium contains 47 expression profiles of natural accessions from diverse geographic origins (Lempe *et al.*, 2005). The phenotypic trait of interest is time to flowering, which is robustly measured by the number of days to bolting or the number of rosette leaves formed before bolting (Amasino, 2010). Timing of flowering according to local climatic conditions is a major determinant of the plants' reproductive success and an important agronomical trait that greatly affects yield. Therefore, an improved knowledge about genes controlling flowering time is of substantial economic value (Craufurd and Wheeler, 2009).

As a second validation of the CStaR method, we compared it with the plain IDA method [ranking as in (1)], Lasso (Tibshirani, 1996), elastic net (Zou and Hastie, 2005) both using linear models (ranking according to absolute values of estimated regression coefficients) and marginal correlation (ranking according to absolute values of marginal correlation to the response) on a publicly available dataset of gene expression profiles in *S. cerevisiae* (Hughes *et al.*, 2000; Supplementary Section S3). This dataset includes both observational and interventional data obtained under similar conditions. Hence, it forms an excellent basis to assess the performance of methods aimed at estimating total causal effects from observational data, as the effects estimated from the observational data can be compared with the effects inferred from the interventional data. These data were used to validate IDA (Maathuis *et al.*, 2010), and we followed the same approach to validate CStaR. In particular, we used the interventional data to infer the total causal effects of the knock-out genes on the remaining genes and defined the top 5% of the effects that were largest in absolute value as the true positives.

3 RESULTS

3.1 Validation for *A. thaliana*

CStaR scores five known regulators of flowering time (*DWF4*, *FLC*, *FRI*, *RPA2B* and *SOC1*; Amasino, 2010; Domagalska *et al.*, 2007; Xia *et al.*, 2006) in its top 25 (Table 1). In particular, *SOC1*, *FRI* and *FLC* are curated flowering time genes in

Table 1. Top 25 findings by CStaR for the *A. thaliana* data

| | Gene | Summary rank | Median effect | Maximum expression | Error (PCER) | Name/annotation |
|----|------------------|--------------|---------------|--------------------|---------------|------------------------|
| 1 | AT2G45660 | 1 | 0.60 | 5.07 | 0.0032 | <i>SOC1</i> |
| 2 | AT4G24010 | 2 | 0.61 | 5.69 | 0.0033 | <i>ATCSLG1</i> |
| 3 | AT1G15520 | 2 | 0.58 | 5.42 | 0.0033 | <i>PDR12</i> |
| 4 | AT3G02920 | 5 | 0.58 | 7.44 | 0.0041 | <i>RPA2B</i> |
| 5 | AT5G43610 | 5 | 0.41 | 4.98 | 0.0069 | <i>ATSUC6</i> |
| 6 | AT4G00650 | 7 | 0.48 | 5.56 | 0.0051 | <i>FRI</i> |
| 7 | AT1G24070 | 8 | 0.57 | 6.13 | 0.0040 | <i>ATCSLA10</i> |
| 8 | AT1G19940 | 9 | 0.53 | 5.13 | 0.0045 | <i>ATGH9B5</i> |
| 9 | AT3G61170 | 9 | 0.51 | 5.12 | 0.0044 | PPR protein |
| 10 | AT1G32375 | 10 | 0.54 | 5.21 | 0.0045 | F-box protein |
| 11 | AT2G15320 | 10 | 0.50 | 5.57 | 0.0047 | LRR protein |
| 12 | AT2G28120 | 10 | 0.49 | 6.45 | 0.0054 | Nodulin protein |
| 13 | AT2G16510 | 13 | 0.50 | 10.7 | 0.0050 | <i>AVAP5</i> |
| 14 | AT3G14630 | 13 | 0.48 | 4.87 | 0.0056 | <i>CYP72A9</i> |
| 15 | AT1G11800 | 15 | 0.51 | 6.97 | 0.0053 | Endonuclease |
| 16 | AT5G44800 | 16 | 0.32 | 6.55 | 0.0079 | <i>CHR4</i> |
| 17 | AT3G50660 | 17 | 0.40 | 7.60 | 0.0078 | <i>DWF4</i> |
| 18 | AT5G10140 | 19 | 0.30 | 10.3 | 0.0085 | <i>FLC</i> |
| 19 | AT1G24110 | 20 | 0.49 | 4.66 | 0.0071 | Peroxidase |
| 20 | AT2G27350 | 20 | 0.48 | 7.06 | 0.0067 | <i>OTLD1</i> |
| 21 | AT1G27030 | 20 | 0.45 | 10.0 | 0.0075 | Unknown protein |
| 22 | AT2G28680 | 22 | 0.46 | 5.23 | 0.0072 | Cupin protein |
| 23 | AT3G16370 | 23 | 0.43 | 12.4 | 0.0099 | Lipase/hydrolase |
| 24 | AT5G25640 | 23 | 0.33 | 5.59 | 0.0091 | Serine protease |
| 25 | AT1G30120 | 24 | 0.46 | 9.97 | 0.0077 | <i>PDH-E1 BETA</i> |

The genes are ranked by increasing summary rank, where ties are sorted according to the estimated median total causal effect taken over 100 stability runs (third column). The maximum expression is taken over the original log₂ data. The error (PCER) is the median PCER over the range of *q* values. *SOC1*, *FRI* and *FLC* are 3 of 119 curated flowering time genes in the Arabidopsis Reactome (Tsesmetzis *et al.*, 2008) (<http://www.arabidopsisreactome.org>). This is a highly significant enrichment of known curated regulators when compared with random guessing ($p < 10^{-5}$, hypergeometric test). Although not curated in Arabidopsis Reactome, also *RPA2B* and *DWF4* are known to affect flowering time (Domagalska *et al.*, 2007; Xia *et al.*, 2006). Since the ordering of the genes in the table is given by their summary rank, the values of median total causal effect and PCER are not decreasing monotonously. For instance, *ATSUC6* has a smaller median total causal effect and a larger PCER than the endonuclease, but since its lower bound for the total causal effect is more stable, the former is ranked 10 positions higher than the latter. All genes from this list, for which mutant seeds were readily available and which were not already known to control flowering time, were used in the subsequent intervention experiments (indicated in bold). In total, intervention experiments were performed for 13 of the 25 top genes not previously known to regulate flowering (Supplementary Section S2.3).

Arabidopsis Reactome (Tsesmetzis *et al.*, 2008) containing 119 known regulators of flowering. This is a highly significant enrichment of known curated regulators when compared with random guessing ($p < 10^{-5}$ in a hypergeometric test). Interestingly, *FLC* and *FRI* are not only major regulators of flowering time in the model species *A. thaliana* but also in the oil-seed rape crop.

Among the other genes in the top 25, which were not already known to play a role in flowering time, there were 13 genes for which mutant seeds were readily available (Supplementary Table S1). These mutants were used for intervention experiments in order to further validate CStaR and to discover new influential genes for flowering time in *A. thaliana* (Supplementary Section S2.3).

The intervention experiments were performed under two photoperiod conditions, short-day (SD) and long-day (LD) with 8 h and 16 h of light, respectively. As phenotypic responses, the number of days to bolting (DTB, for both SD and LD) as well as the rosette leave number (RLN, only for LD) were recorded. Seed viability varied between different genotypes

(Supplementary Tables S2–S4) reducing the number of testable mutants to nine (Supplementary Table S1).

Differences between the knock-out and control group were tested using a two-sided Welch’s *t*-test, because the mutant samples showed different empirical variances compared with the control group. This is most pronounced in the short-day layout. Four new genes were found to have a significant total causal effect on the phenotypic responses at level $\alpha = 0.05$ in at least one of the three settings (Table 2). Among the significant genes is *OTLD1*, a gene involved in chromatin modifications, which may potentially regulate *FLC* expression. Another significant gene is *PDH-E1*, which is involved in carbohydrate metabolism, a known regulation point of flowering time. We did not adjust these *p*-values for multiple testing because we only perform a small number of tests and, in view of small sample sizes, we do not want to sacrifice power. Future studies of the identified novel genes may increase the biological understanding of flowering time control and provide potential targets for breeding strategies in crops. The entire approach from modelling to biological experiments and findings is schematically described in Figure 1.

Table 2. *p*-values from two-sided Welch's *t*-tests in the *A. thaliana* validation

| Gene | Welch's <i>t</i> -test | | |
|---------------------------|------------------------|-------------|-------------|
| | DTB-SD | DTB-LD | RLN-LD |
| <i>PDH-E1 BETA</i> | 0.04 | 0.04 | 0.91 |
| <i>ATGH9B5</i> | 0.02 | 0.15 | 0.04 |
| LRR protein | 0.66 | 0.03 | 0.47 |
| <i>OTLD1</i> | 0.43 | 0.03 | 0.86 |
| <i>PDR12</i> | 0.26 | 0.92 | 0.77 |
| F-box protein | 0.18 | – | – |
| peroxidase | 0.18 | – | – |
| PPR protein | – | 0.65 | 0.47 |
| cupin protein | – | 0.12 | 0.93 |

Only genes are shown for which the insertion was experimentally verified and for which in at least one of the following three settings at least four replicates could be harvested for validation: days to bolting in short days (DTB-SD), days to bolting in long days (DTB-LD) and rosette leave number in long days (RLN-LD). Each mutant was tested versus a control group. *p*-values < 0.05 are written in bold (for complete results see Supplementary Tables S2–S4). A missing entry indicates insufficient number of replicates for testing, i.e. less than four plants.

3.2 Validation for *S. cerevisiae*

We trained the plain IDA method, Lasso (Tibshirani, 1996), elastic net (Zou and Hastie, 2005) and marginal correlation ranking on the observational data, and compared their receiver operating characteristic curves on absolute scale (Fig. 2) showing a clear improvement of CStaR over plain IDA. Moreover, CStaR and IDA are clearly superior to high-dimensional regression methods and marginal correlation screening, which is in line with the earlier validation of IDA (Maathuis *et al.*, 2010).

4 DISCUSSION

We propose CStaR as a general method to obtain a stable ranking of genes in terms of the strengths of their total causal effects on a phenotype of interest. An added value of our method is that, under some assumptions, this ranking comes with an error measure controlling false-positive selections. We showed that CStaR exhibits a large increase in sensitivity when compared with plain IDA and modern regression-type methods in *S. cerevisiae* (Fig. 2). Moreover, we demonstrated the success of CStaR for the biologically much more complex multicellular organism *A. thaliana*. However, in view of uncheckable assumptions (Supplementary Section S1), CStaR is not a tool for confirmatory causal inference.

We used insertion mutant lines for experimental validation. This approach can provide very strong evidence for hypotheses about gene function but it often suffers from a high false-negative rate. Genetic networks are characterized by a high degree of functional redundancy, which can buffer effects of single mutations. The *A. thaliana* genome, for instance, underwent a relatively recent duplication causing partial redundancy between many orthologous gene pairs. Thus, often double mutants need to be tested to observe alterations in phenotype. In addition, the function of essential genes cannot be tested

with insertion mutants. Therefore, the high proportion of confirmation in the test set of insertion mutants is highly reassuring. This makes it plausible that CStaR is relevant for commercial breeding and transgenic approaches. In fact since CStaR is mathematically justified under clearly stated assumptions (Maathuis *et al.*, 2009; Meinshausen and Bühlmann, 2010), it has the potential to generalize many other settings in biology, agriculture and other fields where efficient design and prioritization of new intervention experiments is a core aim.

ACKNOWLEDGEMENTS

We thank T. Wey for the technical assistance with plant work.

Funding: The work was partly financed with a grant of the Swiss SystemsX.ch Initiative to the project 'LiverX' of the Competence Centre for Systems Physiology and Metabolic Diseases. The LiverX project was evaluated by the Swiss National Science Foundation.

Conflict of Interest: none declared.

REFERENCES

- Amasino, R. (2010) Seasonal and developmental timing of flowering. *Plant J.*, **61**, 1001–1013.
- Craufurd, P.Q. and Wheeler, T.R. (2009) Climate change and the flowering time of annual crops. *J. Exp. Bot.*, **60**, 2529–2539.
- Dai, M. *et al.* (2005) Evolving gene/transcript definitions significantly alter the interpretation of GeneChip data. *Nucleic Acids Res.*, **33**, e175.
- Domagalska, M.A. *et al.* (2007) Attenuation of brassinosteroid signaling enhances FLC expression and delays flowering. *Development*, **134**, 2841–2850.
- Gilad, Y. *et al.* (2008) Revealing the architecture of gene regulation: the promise of eQTL studies. *Trends Genet.*, **24**, 408–415.
- Glazier, A.M. *et al.* (2002) Finding genes that underlie complex traits. *Science*, **298**, 2345–2349.
- Hughes, T.R. *et al.* (2000) Functional discovery via a compendium of expression profiles. *Cell*, **102**, 109–126.
- Kalisch, M. *et al.* (2012) Causal inference using graphical models with the R package pcalg. *J. Stat. Softw.*, **47**, 1–26.
- Kliebenstein, D. (2009) Quantitative genomics: analyzing intraspecific variation using global gene expression polymorphisms or eQTLs. *Annu. Rev. Plant Biol.*, **60**, 93–114.
- Kruglyak, L. and Storey, J.D. (2009) Cause and express. *Nat. Biotechnol.*, **27**, 544–545.
- Lempe, J. *et al.* (2005) Diversity of flowering responses in wild *Arabidopsis thaliana* strains. *PLoS Genet.*, **1**, 109–118.
- Maathuis, M.H. *et al.* (2009) Estimating high-dimensional intervention effects from observational data. *Ann. Stat.*, **37**, 3133–3164.
- Maathuis, M.H. *et al.* (2010) Predicting causal effects in large-scale systems from observational data. *Nat. Met.*, **7**, 247–248.
- Meinshausen, N. and Bühlmann, P. (2010) Stability selection. *J. Roy. Stat. Soc. B Met.*, **72**, 417–473.
- Pearl, J. (2000) Causality: models, reasoning and inference, Cambridge Univ. Press, 47.
- Spirtes, P., Glymour, C.N. and Scheines, R. (2000) Causation, prediction and search, The MIT Press, 81.
- Tibshirani, R. (1996) Regression shrinkage and selection via the Lasso. *J. Roy. Stat. Soc. B Met.*, **58**, 267–288.
- Tsesmetzis, N. *et al.* (2008) Arabidopsis reactome: a foundation knowledgebase for plant systems biology. *Plant Cell*, **20**, 1426–1436.
- Xia, R. *et al.* (2006) ROR1/RPA2A, a putative replication protein A2, functions in epigenetic gene silencing and in regulation of meristem development in Arabidopsis. *Plant Cell*, **18**, 85–103.
- Zou, H. and Hastie, T. (2005) Regularization and variable selection via the elastic net. *J. Roy. Stat. Soc. B Met.*, **67**, 301–320.

Supplementary: Causal Stability Ranking

Daniel J. Stekhoven^{1,2,3*}, Izabel Moraes⁴, Gardar Sveinbjörnsson¹,
Lars Hennig^{4,5}, Marloes H. Maathuis¹ and Peter Bühlmann^{1,3}

¹ Seminar for Statistics, Department of Mathematics, ETH Zurich, 8092 Zurich, Switzerland.

² Life Science Zurich PhD Program on Systems Biology of Complex Diseases.

³ Competence Center for Systems Physiology and Metabolic Diseases, 8092 Zurich, Switzerland.

⁴ Uppsala BioCenter, Department of Plant Biology and Forest Genetics, Swedish University of Agricultural Sciences and Linnean Center for Plant Biology, 750 07 Uppsala, Sweden.

⁵ Plant Biotechnology, Department of Biology, ETH Zurich, 8092 Zurich, Switzerland.

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXXX

1 ASSUMPTIONS

There are two main assumptions underlying CStaR.

Gaussian distribution faithful to a DAG and no hidden confounders.

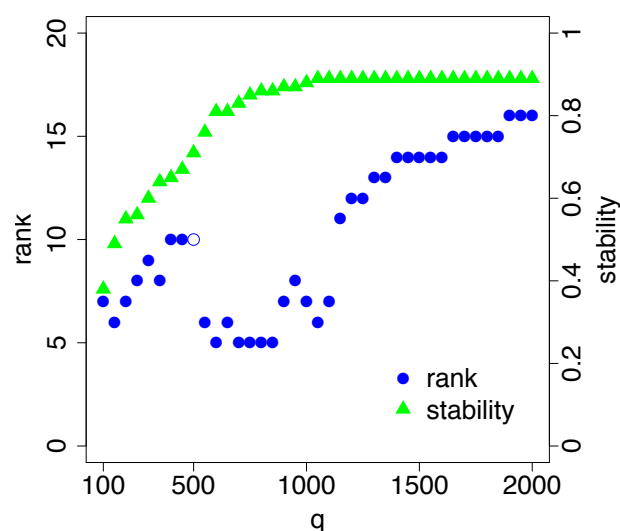
The observational data are assumed to be a sample of n i.i.d. jointly Gaussian random vectors (with dimension $p + 1$; one response and a p -dimensional covariate). Furthermore, this multivariate Gaussian distribution is Markovian and faithful (Spirtes *et al.* (2000)) to the true unknown underlying causal directed acyclic graph (DAG). This implicitly requires us to observe all the relevant variables and that there are no hidden confounders.

The Gaussianity assumption is made in IDA for two reasons: first, it allows us to use partial correlations to test conditional independencies in the PC-algorithm; second, it implies that conditional expectations are linear such that total causal effects can be computed via least squares regression. For details we refer to earlier work (Maathuis *et al.* (2009)). Checking multivariate Gaussianity is merely impossible in high dimensions. However, we verified normality for the marginal distributions using a Lilliefors test with Benjamini-Hochberg multiple testing correction for controlling the false discovery rate, showing approximate normality in the *A. thaliana* example.

The assumption about faithfulness to the true underlying causal DAG and no hidden confounders is more fundamental and is required for the identifiability of the lower bounds for the total causal effects. Without these assumptions we could not make any statements on causal effects from observational data. On the other hand, if we knew the true underlying DAG and there were no hidden confounders we could fully identify the causal effects. However, faithfulness is uncheckable, and there is clearly a limitation that a DAG does not allow for feedback-loops. Also, the assumption of having no unmeasured confounder variables cannot be checked. However, we cover about 80% of all genes in *A. thaliana* and *S. cerevisiae*, and therefore we expect that many important variables are captured¹.

*to whom correspondence should be addressed

¹ We believe that we do not lose major confounding variables when using a correlation pre-screening as used in the validation with *A. thaliana*.



Supplementary Figure 1. Example for the graphical summary of ranks and stabilities of a protein-coding gene (F-box family protein, rank 10). The ranks (blue dots) are assessed for each q -value, as is the stability (green triangles), i.e., the relative frequency $\hat{\Pi}_j$ of this gene j to be scored in the top q . The blue circle at $q = 500$ indicates the median rank.

Exchangeability and “better than random guessing”.

The exchangeability condition requires that the selection of noise variables (i.e. variables in S_{false}) in IDA is equally likely. The “better than random guessing” condition assumes that IDA is performing better than blindly guessing the total causal strength of genes (covariates).

Such a screening reduces the amount of coverage and was done due to computational reasons.

Both assumptions are used for the error control in stability selection with the bound

$$\mathbb{E}[V] \leq \frac{1}{2\pi_{\text{thr}} - 1} \frac{q^2}{p}. \quad (1)$$

For details we refer to earlier work (Meinshausen and Bühlmann (2010)). As with faithfulness, the assumption of exchangeability is uncheckable (inclining on the dependencies of the variables in a very complicated way) but stability selection with the formula in (1) is rather robust with respect to failure of the exchangeability condition (Meinshausen and Bühlmann (2010)). Finally, assuming that IDA is performing better than random guessing seems very plausible (Fig. 2).

While these mathematical conditions are at best only approximately true, we show that CStaR is able to identify the genes (covariates) with a strong total causal influence on a certain phenotype (Table 1, Fig. 2). However, it is merely impossible to validate the PCER measure - controlling a type I statistical error - with real data. While we believe that the numbers for the PCER provide some useful information, they should not be over-interpreted as their validity depends on uncheckable assumptions.

2 VALIDATION FOR *A. THALIANA*

2.1 Description and pre-processing of the data

We use a publicly available compendium of $n = 47$ *Arabidopsis thaliana* gene expression profiles of 4-day old seedlings from a set of 31 wild-type accessions (23 single and 8 triplicate profiles) for which also flowering time data (averaged rosette leave numbers) were available (Lempe *et al.* (2005)), which served as response variable. Microarray raw data, which were based on Affymetrix ATH1 arrays, were downloaded from ArrayExpress (<http://www.ebi.ac.uk/arrayexpress/>, accession E-TABM-18). Expression summaries for 21,440 *A. thaliana* genes were calculated using RMA and the redefined custom CDF ATH1121501_At.TAIRG (Version 12) (Dai *et al.* (2005)). Unique profiles with a maximum \log_2 -expression level above 4.5 were pre-selected, leaving 19,493 genes for the analysis.

2.2 CStaR parameters

Due to computational issues, we reduce the number of genes (covariates) by filtering those which have a marginal correlation of at least 0.4 with the response of interest in each stability run. Furthermore, the data were standardized such that the expressions for each gene (over the different samples) have mean zero and standard deviation one in each stability run. The tuning parameter for the PC-algorithm in IDA is set to $\alpha = 0.1$. The summary rank is aggregated over $q \in \{100, 150, 200, \dots, 2,000\}$. The upper limit of 2,000 is chosen because this is about the maximum dimension after filtering according to marginal correlations in each subsample.

2.3 Biological validation

Seeds of Columbia (Col) *Arabidopsis thaliana* wild-type and T-DNA insertion mutants were obtained from the Nottingham Arabidopsis Stock Centre (www.arabidopsis.info). Homozygous insertion mutants were selected for top-scoring candidate

Supplementary Table 1. List of mutant lines used for biological validation. The second column gives the summary rank from Table 1 and the fourth column the name of the homozygous insertion mutant line (Supplementary Section 2.3). Mutant lines for which in both experimental layout (SD and LD) not enough viable replicates (less than 4) could have been collected are indicated with a dagger (†). In case of *ATSUC6* and the nodulin protein (‡) it was not possible to confirm the annotated knock out.

| gene | summary rank | name | mutant line |
|------------------|--------------|--------------------|---------------|
| <i>AT1G15520</i> | 2 | <i>PDR12</i> | SALK_013945Cc |
| <i>AT5G43610</i> | 5 | <i>ATSUC6</i> ‡ | SALK_132450c |
| <i>AT1G24070</i> | 8 | <i>ATCSLA10</i> † | SALK_023438Cc |
| <i>AT1G19940</i> | 9 | <i>ATGH9B5</i> | GK-280G04.9c |
| <i>AT3G61170</i> | 9 | PPR protein | SALK_013940Cc |
| <i>AT2G15320</i> | 10 | LRR protein | SALK_032550Cc |
| <i>AT1G32375</i> | 10 | F-box protein | SALK_096159Cc |
| <i>AT2G28120</i> | 10 | nodulin protein ‡ | SK36217 |
| <i>AT1G11800</i> | 15 | nuclease † | SALK_043413Cc |
| <i>AT2G27350</i> | 20 | <i>OTLD1</i> | SALK_037047Cc |
| <i>AT1G24110</i> | 20 | peroxidase | SALK_087392Cc |
| <i>AT2G28680</i> | 22 | cupin protein | SALK_074581Cc |
| <i>AT1G30120</i> | 24 | <i>PDH-E1 BETA</i> | SALK_046011Cc |

genes that were not previously known to affect flowering time. For 13 such mutant lines a sufficient number of viable seeds for experimental confirmation was obtained (see Supplementary Table 1). For measuring flowering time, seeds were plated on Murashige and Skoog (MS) medium (Duchefa, Haarlem, The Netherlands), stratified for 2 days at 4°C, and grown on plates for 10 days before transfer onto soil. Plants were kept in Conviron growth chambers with mixed cold fluorescent and incandescent light (110 to 140 $\mu\text{mol/m}^2\text{s}$, $21 \pm 2^\circ\text{C}$) under long day (LD, 16h) or short day (SD, 8h) photoperiods. The flowering time was measured as days to bolting (DTB, LD and SD) and as the number of total rosette leaves at bolting (RLN, LD only). Rosette leaves were not counted for the SD experiment because the plants develop up to 100 leaves under such conditions. Because of low numbers of viable seeds, the number of testable mutant lines got reduced to nine (i.e. 4 or more plants were required for subsequent testing). Furthermore, DTB of some lines could only be tested under LD or SD (for full details see also Supplementary Tables 2, 3 and 4).

2.4 Analysis of validation results in *A. thaliana*

We tested the mutant groups versus the wild-type control group using a two-sided Welch's t-test in all three cases (i.e. SD in days, Supplementary Table 2; LD in days, Supplementary Table 3; and LD in RLN, Supplementary Table 4). The Welch's t-test assumes different empirical variances in the groups which we consider to be reasonable.

Supplementary Table 2. Short-day layout (DTB-SD). Results from two-sided Welch's t-test considering the flowering time. The mean duration in the control group was 87.3 days coming from 6 replicates. The columns give gene names, mean number of days to bolting, t-statistic, degrees of freedom, number of replicates and the p-value. Mutant lines for which in both experimental layout (SD and LD) not enough viable replicates (less than 4) could have been collected are indicated with a dagger (†).

| annotation | mean | t | df | #rep | p-value |
|----------------------------------|------------------------|------|------|------|---------|
| <i>PDR12</i> | 82.7 | 1.22 | 7.90 | 7 | 0.26 |
| <i>ATGH9B5</i> | 76.9 | 2.83 | 7.51 | 14 | 0.02 |
| F-box protein | 81.8 | 1.49 | 7.25 | 4 | 0.18 |
| <i>OTLD1</i> | 84.1 | 0.84 | 8.56 | 9 | 0.43 |
| peroxidase | 82.0 | 1.51 | 6.38 | 8 | 0.18 |
| <i>PDH-E1 BETA</i> | 78.0 | 2.64 | 6.48 | 7 | 0.04 |
| <i>ATSUC6</i> nodulin protein | mutation not confirmed | | | | |
| <i>ATCSLA10</i> † | 86.5 | - | - | 2 | - |
| PPR protein† | 83.0 | - | - | 1 | - |
| LRR protein | 84.7 | - | - | 3 | - |
| cupin protein† | 84.0 | - | - | 2 | - |
| nuclease† | 76.0 | - | - | 1 | - |

Supplementary Table 3. Long-day layout (DTB-LD). Results from two-sided Welch's t-test considering the flowering time. The mean duration in the control group was 26.1 days coming from 11 replicates. The columns give gene names, mean number of days to bolting, t-statistic, degrees of freedom, number of replicates and the p-value. Mutant lines for which in both experimental layout (SD and LD) not enough viable replicates (less than 4) could have been collected are indicated with a dagger (†).

| annotation | mean | t | df | #rep | p-value |
|----------------------------------|------------------------|-------|-------|------|---------|
| <i>PDR12</i> | 26.0 | 0.10 | 10.35 | 7 | 0.92 |
| <i>ATGH9B5</i> | 24.8 | 1.58 | 9.88 | 6 | 0.15 |
| PPR protein | 25.6 | 0.46 | 12.96 | 10 | 0.65 |
| LRR protein | 24.1 | 2.43 | 18.77 | 12 | 0.03 |
| <i>OTLD1</i> | 24.4 | 2.24 | 23.75 | 15 | 0.03 |
| cupin protein | 27.4 | -1.74 | 8.79 | 5 | 0.12 |
| <i>ATSUC6</i> nodulin protein | mutation not confirmed | | | | |
| <i>ATCSLA10</i> † | 28.0 | - | - | 1 | - |
| F-box protein† | - | - | - | 0 | - |
| nuclease† | 25.0 | - | - | 2 | - |
| peroxidase† | 28.0 | - | - | 1 | - |
| <i>PDH-E1 BETA</i> † | 25.0 | - | - | 3 | - |

3 VALIDATION FOR *S. CEREVISIAE*

3.1 Description of the data

We validated CStAR on the same data (Hughes *et al.* (2000)) as in the original IDA article (Maathuis *et al.* (2010)), using the same pre-processing steps. These data contain 234 interventional and 63 observational full-genome expression profiles ($p = 5,361$) in *S. cerevisiae* (*Saccharomyces cerevisiae*). Genes were not further pre-selected using any expression cutoff.

3.2 CStAR parameters

The tuning parameter for the PC-algorithm in IDA was set to $\alpha = 0.01$. The summary rank is aggregated over different proportions of all possible total causal effects of the 234 knock-out genes on the remaining genes (i.e., $5,361 \times 234 - 234 = 1,254,240$). We choose the range of q values in terms of percentages of the number of all possible effects, that is $q \in \{0.01\%, 0.03\%, 0.05\%, \dots, 1\%\} \times 1,254,240$.

3.3 Analysis of validation results in *S. cerevisiae*

We used CStAR to obtain a stable ranking of the total causal effects of the knock-out genes on the remaining genes based on the observational data. As comparison, we also applied IDA, Lasso (Tibshirani (1996)) and elastic net (Zou and Hastie (2005)), the latter two using a linear model, to the observational data, as described in the IDA validation (Maathuis *et al.* (2010)). Moreover, we added marginal correlation screening as an extra competitor and applied this as follows. The correlation $\rho_{i,j}$ of each knock-out gene i ($i = 1, \dots, 234$) with all other genes j ($j = 1, \dots, 5,360$) are

computed. Absolute correlations are then sorted in decreasing order to obtain a ranking.

We used the interventional data to infer the total causal effects of the knock-out genes on the remaining genes (Maathuis *et al.* (2010)). The top 5% of the largest effects in absolute value were defined as the target set of effects that we want to be able to identify. We then compared the receiver operating characteristic (ROC) curves of the different methods on the absolute scale (Fig. 2).

REFERENCES

- Amasino, R. (2010). Seasonal and developmental timing of flowering. *The Plant Journal*, **61**(6), 1001–1013.
- Craufurd, P. and Wheeler, T. (2009). Climate change and the flowering time of annual crops. *Journal of Experimental Botany*, **60**(9), 2529–2539.
- Dai, M., Wang, P., Boyd, A., Kostov, G., Atthey, B., Jones, E., Bunney, W., Myers, R., Speed, T., Akil, H., *et al.* (2005). Evolving gene/transcript definitions significantly alter the interpretation of GeneChip data. *Nucleic Acids Research*, **33**(20), e175.
- Domagalska, M., Schomburg, F., Amasino, R., Vierstra, R., Nagy, F., and Davis, S. (2007). Attenuation of brassinosteroid signaling enhances FLC expression and delays flowering. *Development*, **134**(15), 2841–2850.
- Gilad, Y., Rifkin, S., and Pritchard, J. (2008). Revealing the architecture of gene regulation: the promise of eQTL studies. *Trends in genetics*, **24**(8), 408–415.
- Glazier, A., Nadeau, J., and Aitman, T. (2002). Finding genes that underlie complex traits. *Science*, **298**(5602), 2345–2349.
- Hughes, T., Marton, M., Jones, A., Roberts, C., Stoughton, R., Armour, C., Bennett, H., Coffey, E., Dai, H., He, Y., *et al.* (2000). Functional discovery via a compendium of expression profiles. *Cell*, **102**(1), 109–126.
- Kalisch, M., Mächler, M., Colombo, D., Maathuis, M., and Bühlmann, P. (2012). Causal inference using graphical models with the R package pcalg. *Journal for Statistical Software*, **47**(11), 1–26.
- Kliebenstein, D. (2009). Quantitative genomics: analyzing intraspecific variation using global gene expression polymorphisms or eQTLs. *Plant Biology*, **60**(1), 93.

Supplementary Table 4. Long-day layout (RLN-LD). Results from two-sided Welch's t-test considering the number of leaves. The mean number of rosette leaves was 8.6 coming from 11 replicates. The columns give gene names, mean number of rosette leaves, t-statistic, degrees of freedom, number of replicates and the p-value. Mutant lines for which in both experimental layout (SD and LD) not enough viable replicates (less than 4) could have been collected are indicated with a dagger (†).

| annotation | mean | t | df | #rep | p-value |
|----------------------------------|------------------------|-------|-------|------|---------|
| <i>PDR12</i> | 8.7 | -0.30 | 13.78 | 7 | 0.77 |
| <i>ATGH9B5</i> | 10.0 | -2.30 | 10.02 | 6 | 0.04 |
| PPR protein | 8.9 | -0.73 | 18.82 | 10 | 0.47 |
| LRR protein | 8.2 | 0.73 | 20.95 | 12 | 0.47 |
| <i>OTLD1</i> | 8.5 | 0.19 | 16.69 | 15 | 0.86 |
| cupin protein | 8.6 | -0.09 | 8.30 | 5 | 0.93 |
| <i>ATSUC6</i> nodulin protein | mutation not confirmed | | | | |
| <i>ATCSLA10</i> † | 8.0 | - | - | 1 | - |
| F-box protein † | - | - | - | 0 | - |
| nuclease † | 8.5 | - | - | 2 | - |
| peroxidase † | 9.0 | - | - | 1 | - |
| PDH-E1 BETA † | 8.7 | - | - | 3 | - |

- Kruglyak, L. and Storey, J. (2009). Cause and express. *Nature Biotechnology*, **27**(6), 544–545.
- Lempe, J., Balasubramanian, S., Sureshkumar, S., Singh, A., Schmid, M., and Weigel, D. (2005). Diversity of flowering responses in wild *Arabidopsis thaliana* strains. *PLoS Genetics*, **1**(1), 109–118.
- Maathuis, M., Kalisch, M., and Bühlmann, P. (2009). Estimating high-dimensional intervention effects from observational data. *The Annals of Statistics*, **37**(6A), 3133–3164.
- Maathuis, M., Colombo, D., Kalisch, M., and Bühlmann, P. (2010). Predicting causal effects in large-scale systems from observational data. *Nature Methods*, **7**(4), 247–248.
- Meinshausen, N. and Bühlmann, P. (2010). Stability selection. *Journal of the Royal Statistical Society: Series B*, **72**(4), 417–473.
- Spirites, P., Glymour, C., and Scheines, R. (2000). *Causation, Prediction, and Search*. MIT Press, second edition.
- Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society: Series B*, **58**(1), 267–288.
- Tsesmetzis, N., Couchman, M., Higgins, J., Smith, A., Doonan, J., Seifert, G., Schmidt, E., Vastrik, I., Birney, E., Wu, G., et al. (2008). Arabidopsis reactome: a foundation knowledgebase for plant systems biology. *The Plant Cell Online*, **20**(6), 1426–1436.
- Xia, R., Wang, J., Liu, C., Wang, Y., Wang, Y., Zhai, J., Liu, J., Hong, X., Cao, X., Zhu, J., et al. (2006). ROR1/RPA2A, a putative replication protein A2, functions in epigenetic gene silencing and in regulation of meristem development in *Arabidopsis*. *The Plant Cell Online*, **18**(1), 85–103.
- Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B*, **67**(2), 301–320.