

False Discovery Rate Control for Fast Screening of Large-Scale Genomics Biobanks

Jasin Machkour

Technische Universität Darmstadt
64283 Darmstadt, Germany
jasin.machkour@tu-darmstadt.de

Michael Muma

Technische Universität Darmstadt
64283 Darmstadt, Germany
michael.muma@tu-darmstadt.de

Daniel P. Palomar

The Hong Kong University of Science and Technology
Clear Water Bay, Hong Kong SAR, China
palomar@ust.hk

Abstract—Genomics biobanks are information treasure troves with thousands of phenotypes (e.g., diseases, traits) and millions of single nucleotide polymorphisms (SNPs). The development of methodologies that provide reproducible discoveries is essential for the understanding of complex diseases and precision drug development. Without statistical reproducibility guarantees, valuable efforts are spent on researching false positives. Therefore, scalable multivariate and high-dimensional false discovery rate (FDR)-controlling variable selection methods are urgently needed, especially, for complex polygenic diseases and traits. In this work, we propose the Screen-T-Rex selector, a fast FDR-controlling method based on the recently developed T-Rex selector. The method is tailored to screening large-scale biobanks and it does not require choosing additional parameters (sparsity parameter, target FDR level, etc). Numerical simulations and a real-world HIV-1 drug resistance example demonstrate that the performance of the Screen-T-Rex selector is superior, and its computation time is multiple orders of magnitude lower compared to current benchmark knockoff methods.

Index Terms—Screen-T-Rex selector, FDR control, high-dimensional variable selection, GWAS, HIV-1 drug resistance.

I. INTRODUCTION

The systematic screening of large-scale genomics biobanks enables understanding complex diseases and aids in drug development [1]. Achieving these goals requires finding the few reproducible associations among potentially millions of single nucleotide polymorphisms (SNPs) and a phenotype, i.e., a disease or trait of interest. This allows to further study potentially functionally associated regions on the genome. Large biobanks, such as the UK biobank [2], contain thousands of phenotypes and large ultra-high-dimensional genomics data, where the number of variables (i.e., SNPs) p is much larger than the number of observations n . The above described genome-wide association studies (GWAS) [3] require time and cost intensive follow-up investigations. Hence, it is of utmost importance to keep the number of false discoveries low while discovering as many associations as possible. Therefore, we will consider two metrics:

1. The false discovery rate (FDR) is the expected value of the false discovery proportion (FDP), i.e., the expected percentage of false discoveries among all discoveries: $\text{FDR} := \mathbb{E}[\text{FDP}] := \mathbb{E}[\# \text{ False discoveries} / \# \text{ Discoveries}]$.
2. The true positive rate (TPR) is the expected value of the true positive proportion (TPP), i.e., the expected percentage of true discoveries among all true active variables: $\text{TPR} := \mathbb{E}[\text{TPP}] := \mathbb{E}[\# \text{ True discoveries} / \# \text{ True actives}]$.

Existing FDR-controlling methods allow the user to set a target FDR $\alpha \in [0, 1]$ and select variables such that the FDR is controlled at the target level (i.e., α is not exceeded) while maximizing the number of selected variables and, thus, implicitly maximizing the

The first and second author are supported by the LOEWE initiative (Hesse, Germany) within the emergenCITY center. The second author is also supported by the ERC Starting Grant ScReeningData. The third author is supported by the Hong Kong GRF 16207820 research grant.

Extensive calculations on the Lichtenberg High-Performance Computer of the Technische Universität Darmstadt were conducted for this research.

TPR. Popular methods for low-dimensional settings (i.e., $n \geq p$) are the Benjamini-Hochberg (BH) method [4], the Benjamini-Yekutieli (BY) method [5], and the more recent *fixed-X* knockoff methods [6]. Unfortunately, these methods are not applicable for the multivariate analysis of high-dimensional ($p > n$) settings such as GWAS.

In recent years, multivariate FDR-controlling methods for high-dimensional multivariate GWAS have been proposed: *model-X* (and related) knockoff methods [7]–[10] and *T-Rex* selector methods [11]–[14]. However, Figure 1 in [11] shows that only the *T-Rex* selector is scalable to millions of variables in a reasonable computation time, while the computation time of the *model-X* methods is multiple orders of magnitude higher and, thus, the method becomes practically infeasible in large-scale settings. Nevertheless, even the comparably low computation time of the *T-Rex* selector for one phenotype might become a burden when conducting GWAS for many phenotypes.

Therefore, we propose the *Screen-T-Rex* selector, a fast version of the *T-Rex* selector. The proposed FDR-controlling method is suitable for conducting large-scale GWAS (with up to millions of SNPs) for tens of thousands of phenotypes. It does not ask the user to set a target FDR level, but provides the user with an estimate of the achieved FDR. In the cases, where the user is not satisfied with the provided FDR estimate, the original *T-Rex* selector should be used with the result of the Screen-T-Rex selector and the desired target FDR as inputs. The proposed *Screen-T-Rex* selector has the following three major innovations/advantages:

1. It provably controls the FDR at the self-estimated level (see Theorems 1 and 2 in Section III).
2. It does not require the choice of any additional parameters (sparsity parameter, target FDR level, etc.).
3. Its computation time is approximately one order of magnitude lower than that of the original *T-Rex* selector and more than three orders of magnitude lower than that of the *model-X* knockoff methods in our simulations (see Table I).

Organization: Section II briefly revisits the original *T-Rex* selector. In Section III, the *Screen-T-Rex* selector is proposed. Sections IV, V, and VI, compare the proposed method against benchmark methods via numerical simulations, a simulated GWAS, and a real world HIV-1 drug resistance study, respectively. Section VII concludes the paper.

II. THE T-REX SELECTOR

The *T-Rex* selector [11] is a fast FDR-controlling variable selection method for high-dimensional ($p > n$) as well as low-dimensional ($n \geq p$) settings. Figure 1 shows a simplified sketch of the *T-Rex* selector framework. It requires the following inputs:

1. A predictor matrix $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_p]$ that, e.g., contains p SNPs $\mathbf{x}_1, \dots, \mathbf{x}_p$ as columns, where $\mathbf{x}_j = [x_{1j} \cdots x_{nj}]^T$ contains n observations of the j th SNP. That is, the i th row of \mathbf{X} contains the measurements of all SNPs for the i th subject.
2. A response vector $\mathbf{y} = [y_1 \cdots y_n]^T$ that, e.g., contains the phenotypes of all n subjects. These can be measurements of

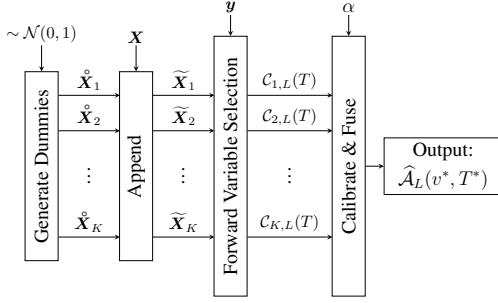


Figure 1: Sketch of the T -Rex selector framework [12].

the disease progression or, in a simple case-control study, the value “1” for cases and the value “0” for controls.

3. The target FDR level $\alpha \in [0, 1]$.

First, the T -Rex selector generates K dummy predictor matrices $\hat{\mathbf{X}}_k = [\hat{x}_1 \cdots \hat{x}_L]$, $k = 1, \dots, K$ (containing L dummy predictors), where each element of the l th dummy predictor $\hat{x}_l = [\hat{x}_{1l} \cdots \hat{x}_{nl}]^\top$ is sampled from a univariate standard normal distribution. After appending each dummy matrix to the original predictor matrix, K independent random experiments are conducted by applying a forward variable selection algorithm to each extended predictor matrix $\tilde{\mathbf{X}}_k = [\mathbf{X} \ \hat{\mathbf{X}}_k]$, $k = 1, \dots, K$, where the response \mathbf{y} acts as the supervising vector. The forward variable selection algorithm includes one variable at a time and terminates after $T \geq 1$ dummies have been included. In [11], it is proposed to use the LARS algorithm [15] (or related methods, see, e.g., [16]–[18]), which assumes a linear relationship between the predictors and the response, as the forward selector within the T -Rex framework. Following the notation of the T -Rex selector, the linear model is defined by

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (1)$$

where $\boldsymbol{\beta}$ is the sparse coefficient vector and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ is the Gaussian noise vector. The obtained candidate sets $\mathcal{C}_{k,L}(T)$, $k = 1, \dots, K$, contain the included original variables (after removal of the T dummy variables). In the “Calibrate & Fuse” step, the relative occurrence of each original variable j in the candidate sets, i.e., $\Phi_{T,L}(j) \in [0, 1]$, $j = 1, \dots, p$, is computed. The T -Rex calibration algorithm automatically determines

1. the number of dummies L ,
2. the optimal number of included dummies before terminating the random experiments T^* , and
3. the optimal voting level $v^* \in [0, 1]$

to obtain the set of selected variables

$$\hat{\mathcal{A}}_L(v^*, T^*) := \{j : \Phi_{T^*,L}(j) > v^*\}$$

such that the FDR is controlled at the target level α , while the number of selected variables $|\hat{\mathcal{A}}_L(v^*, T^*)|$ is maximized. Note that in accordance with the suggestion in [11], we have conducted $K = 20$ random experiments throughout this work.

Definition 1 (FDP and FDR). Let $R_{T,L}(v) := |\hat{\mathcal{A}}_L(v, T)|$ and $V_{T,L}(v) := |\{\text{null } j : \Phi_{T,L}(j) > v\}|$ be the number of selected variables and the number of selected null variables (i.e., false positives), respectively. Define $a \vee b := \max\{a, b\}$, $a, b \in \mathbb{R}$. Then, the FDR and the false discovery proportion (FDP) are defined by

$$\text{FDR}(v, T, L) := \mathbb{E}[\text{FDP}(v, T, L)] := \mathbb{E}\left[\frac{V_{T,L}(v)}{R_{T,L}(v) \vee 1}\right].$$

III. PROPOSED: SCREEN-T-REX SELECTOR

Two versions of the *Screen-T-Rex* selector are proposed, the corresponding FDR control theorems are presented, and an algorithm for screening genomics biobanks is formulated.

A. Ordinary Screen-T-Rex Selector

While the original T -Rex selector determines T , L , and v such that the FDR is controlled at the user-defined target level, the *Screen-T-Rex* selector fixes $(T, L, v) = (1, p, 0.5)$. This is a special case of the original T -Rex selector that

1. is harnessed by the proposed *Screen-T-Rex* selector to determine an estimator of the FDR and
2. requires a much lower computation time than the original T -Rex selector and other benchmark methods (see Table I in Section V).

The FDR estimator of the proposed ordinary *Screen-T-Rex* selector is given by

$$\hat{\alpha} := 1/(R_{1,p}(0.5) \vee 1),$$

i.e., one divided by the number of selected variables. The intuition behind this estimator is as follows: $T = 1$ dummy variable is allowed to enter the solution paths of the random experiments before terminating the forward selection processes. So, in each random experiment one out of p dummies is included. Therefore, we expect, on average, no more than one out of at most p null variables to be included in each candidate set $\mathcal{C}_{k,L}(T)$, and, consequently, no more than one null variable among all selected variables. This idea is formalized in the following FDR control result:

Theorem 1 (FDR control - ordinary *Screen-T-Rex*). Define $\hat{\alpha} := 1/(R_{1,p}(0.5) \vee 1)$. Then, $\text{FDR} = \mathbb{E}[\text{FDP}] \leq \hat{\alpha}$, i.e., the FDR is controlled at the estimated level $\hat{\alpha}$.¹

Proof. With Definition 1, we obtain

$$\text{FDP} = \frac{V_{1,p}(0.5)}{R_{1,p}(0.5) \vee 1} = \hat{\alpha} \cdot V_{1,p}(0.5). \quad (2)$$

Let $p = p_1 + p_0$, where p_1 and p_0 are the number of true active and null variables, respectively. Taking the expectation of (2) yields

$$\text{FDR} = \mathbb{E}[\text{FDP}] = \hat{\alpha} \cdot \mathbb{E}[V_{1,p}(0.5)] \leq \hat{\alpha} \cdot \frac{p_0}{p+1} \leq \hat{\alpha},$$

where the first inequality follows from $V_{1,p}(0.5)$ being stochastically dominated by the negative hypergeometric distribution $\text{NHG}(p_0 + p, p_0, 1)$, whose expected value is given by $p_0/(p+1)$ (for details on the NHG, see [11]). \square

B. Confidence-Based Screen-T-Rex Selector

The above proposed ordinary *Screen-T-Rex* selector, as well as the original T -Rex selector, only considers the relative occurrences of the candidate variables in the selected active sets $\mathcal{C}_{k,L}(T)$ and disregards the original and dummy coefficient estimates, i.e.,

1. $\hat{\beta}_{j,k}(T, L)$, $j = 1, \dots, p$, (i.e., coefficient estimate of the j th original variable in the k th random experiment and
2. $\hat{\beta}_{l,k}(T, L)$, $l = 1, \dots, L$, (i.e., coefficient estimate of the l th dummy variable in the k th random experiment.

However, since the dummy variables act as flagged null variables (for details, see [11]), the coefficients of the dummies contain information about the distribution of the coefficients of the null variables. Therefore, we propose to harness the coefficient estimates of the dummies to construct a confidence interval

$$C(\gamma) := [c_1(\gamma), c_2(\gamma)], \quad \gamma \in [0, 1], \quad (3)$$

¹For simplicity, $\text{FDR} := \text{FDR}(1, p, 0.5)$ and $\text{FDP} := \text{FDP}(1, p, 0.5)$.

where $c_1(\gamma)$ and $c_2(\gamma)$ are the lower and upper bound, respectively, and γ is the confidence level. The coefficient estimates of the null variables can also be expected to lie within the same confidence interval. Therefore, instead of selecting variables based on their relative occurrences, we replace $V_{1,p}(0.5)$ and $R_{1,p}(0.5)$ in Definition 1 and Theorem 1 by

$$V_{1,p}^{(C)}(\gamma) := |\{\text{null } j : \bar{\beta}_j(1,p) \notin C(\gamma)\}| \text{ and} \quad (4)$$

$$R_{1,p}^{(C)}(\gamma) := |\widehat{\mathcal{A}}_p^{(C)}(\gamma, 1)| := |\{j : \bar{\beta}_j(1,p) \notin C(\gamma)\}|, \quad (5)$$

respectively, where $\bar{\beta}_j(1,p) := \frac{1}{K} \sum_{k=1}^K \hat{\beta}_{j,k}(1,p)$. That is, only candidate variables whose averaged (over K random experiments) coefficient estimates are not inside the confidence interval $C(\gamma)$ are selected.

We propose to construct the confidence interval in (3) using the non-parametric bootstrap with 1,000 resamples of the vector containing the $K = 20$ non-zero dummy coefficient estimates. Since, in all our simulations, the distribution of the bootstrapped standard errors of the averaged non-zero dummy coefficient estimates followed the standard normal distribution, we construct a normal bootstrap confidence interval (for details, see [19], [20]). In the following theorem, we state how the most liberal confidence level γ can be determined such that the FDR is controlled at the estimated level by the confidence-based *Screen-T-Rex* selector:

Theorem 2 (FDR control - confidence-based *Screen-T-Rex*). *Define* $\gamma := \inf \{\gamma' \in [0, 1] : R_{1,p}^{(C)}(\gamma') \leq R_{1,p}(0.5)\}$ and $\hat{\alpha}_C := 1/(R_{1,p}^{(C)}(\gamma) \vee 1)$. *Suppose that* $V_{1,p}^{(C)}(\gamma) \stackrel{d}{\leq} V_{1,p}(0.5)$, *where* $\stackrel{d}{\leq}$ *denotes stochastic dominance. Then,* $\text{FDR} = \mathbb{E}[\text{FDP}] \leq \hat{\alpha}_C$.

Proof. With Definition 1 and Equations (4) and (5), we obtain

$$\begin{aligned} \text{FDR} = \mathbb{E}[\text{FDP}] &= \mathbb{E} \left[\frac{V_{1,p}^{(C)}(\gamma)}{R_{1,p}^{(C)}(\gamma) \vee 1} \right] = \hat{\alpha}_C \cdot \mathbb{E} [V_{1,p}^{(C)}(\gamma)] \\ &\leq \hat{\alpha}_C \cdot \mathbb{E} [V_{1,p}(0.5)] \leq \hat{\alpha}_C, \end{aligned}$$

where the first inequality follows from $V_{1,p}^{(C)}(\gamma) \stackrel{d}{\leq} V_{1,p}(0.5)$ and the second inequality is the same as in the proof of Theorem 1. \square

C. Screening Genomics Biobanks

The *Screen-T-Rex* selector is intended to be used for screening thousands of phenotypes in large biobanks, while only using the original *T-Rex* selector in the cases where the estimated FDR is not acceptable to the user. Here, the user sets the target FDR for the original *T-Rex* selector and a lower and upper bound α_l and α_u , respectively, for the estimated FDRs by both versions of the *Screen-T-Rex* selector. The lower bound is required to avoid solutions at very

low estimated FDRs, since these would yield a low power (i.e., TPR). Algorithm 1 summarizes the proposed work flow.

IV. NUMERICAL EXPERIMENTS

We simulate a high-dimensional data setting according to the linear model in (1) with $n = 300$ samples and $p = 1,000$ predictors (i.e., candidate variables), and $p_1 = 10$ true active variables. The noise variance σ^2 is chosen such that the signal-to-noise-ratio $\text{SNR} := \text{Var}(\mathbf{X}\boldsymbol{\beta}) / \text{Var}(\boldsymbol{\epsilon})$ (for details, see [13]) takes on the values on the x -axes in Figure 2. Note that the FDP and TPP in Figure 2 are averaged over 955 Monte Carlo replications, respectively, and, therefore, are estimates of the FDR and TPR, respectively. A discussion of the results is provided within the caption of Figure 2.

Algorithm 1 Screening Genomics Biobanks.

1. **Input:** α , α_l , and α_u .
2. For each considered phenotype in the biobank do:
 - 2.1. **Run** the *Screen-T-Rex* selector and obtain the estimated FDR levels $\hat{\alpha}$ and $\hat{\alpha}_C$.
 - 2.2. **Determine** the final set of selected variables $\hat{\mathcal{A}}$ as follows:

$$\hat{\mathcal{A}} := \begin{cases} \hat{\mathcal{A}}_p^{(C)}(\gamma, 1), & \alpha_l \leq \hat{\alpha}_C \leq \alpha_u \ \& \\ & \max\{\hat{\alpha}_C, \hat{\alpha} \cdot I(\hat{\alpha} \leq \alpha_u)\} = \hat{\alpha}_C \\ \hat{\mathcal{A}}_p(0.5, 1), & \alpha_l \leq \hat{\alpha} \leq \alpha_u \ \& \\ \emptyset, & \max\{\hat{\alpha}_C \cdot I(\hat{\alpha}_C \leq \alpha_u), \hat{\alpha}\} = \hat{\alpha} \\ & \text{otherwise} \end{cases},$$

where $I(a \leq b)$, $a, b \in \mathbb{R}$, is the indicator function that has the value one if $a \leq b$ and zero otherwise. \emptyset denotes the empty set. Convention: If $\hat{\alpha} = \hat{\alpha}_C$ and all conditions in the first two cases are satisfied, then $\hat{\mathcal{A}} := \hat{\mathcal{A}}_p^{(C)}(\gamma, 1)$.

- 2.3. If $\hat{\mathcal{A}} = \emptyset$, **run** the *T-Rex* selector with target FDR α and **determine**

$$\hat{\mathcal{A}} := \hat{\mathcal{A}}_L(v^*, T^*).$$

3. **Output:** Selected active set $\hat{\mathcal{A}}$.

V. SIMULATED GWAS

The simulation setup and preprocessing of the data in this section is the same as the one in [11]. That is, 100 genomics data sets are simulated using the software HAPGEN2 [21]. It takes haplotypes from the HapMap 3 project [22] as an input and generates the predictor matrix \mathbf{X} that contains $p = 20,000$ SNPs/columns and $n = 1,000$ observations/rows (i.e., 700 cases and 300 controls). The phenotype vector \mathbf{y} contains ones for cases and zeros for controls. Each of the 100 simulated data sets contains $p_1 = 10$ true active

Table I For all methods, the average achieved FDP is lower than the average estimated/target FDR, i.e., all methods control the FDR. The sequential computation time (of both versions) of the proposed *Screen-T-Rex* selector is more than three orders of magnitude lower than that of the *model-X* knockoff+ method. Nearly one order of magnitude is gained compared to the original *T-Rex* selector. Applying Algorithm 1 with $\alpha = 10\%$, $\alpha_l = 5\%$, and $\alpha_u = 20\%$, yields an average FDP and TPP of 15.96% and 47.2%, respectively, without requiring Step 2.3.

Methods	FDR control?	Av. FDP (in %)	Av. estimated/target FDR (in %)	Av. TPP (in %)	Av. sequential comp. time (hh:mm:ss)	Av. relative sequential comp. time
Proposed:						
1. <i>Screen-T-Rex</i> (ordinary)	✓	15.96	18.57	47.2	00:00:44	1
2. <i>Screen-T-Rex</i> (conf.-based)	✓	10.16	12.5	31.7	00:00:45	1.02
Benchmarks:						
3. <i>T-Rex</i>	✓	6.45	10	38.5	00:06:39	8.88
4. <i>model-X+</i>	✓	0	10	0	20:00:38	1601.39

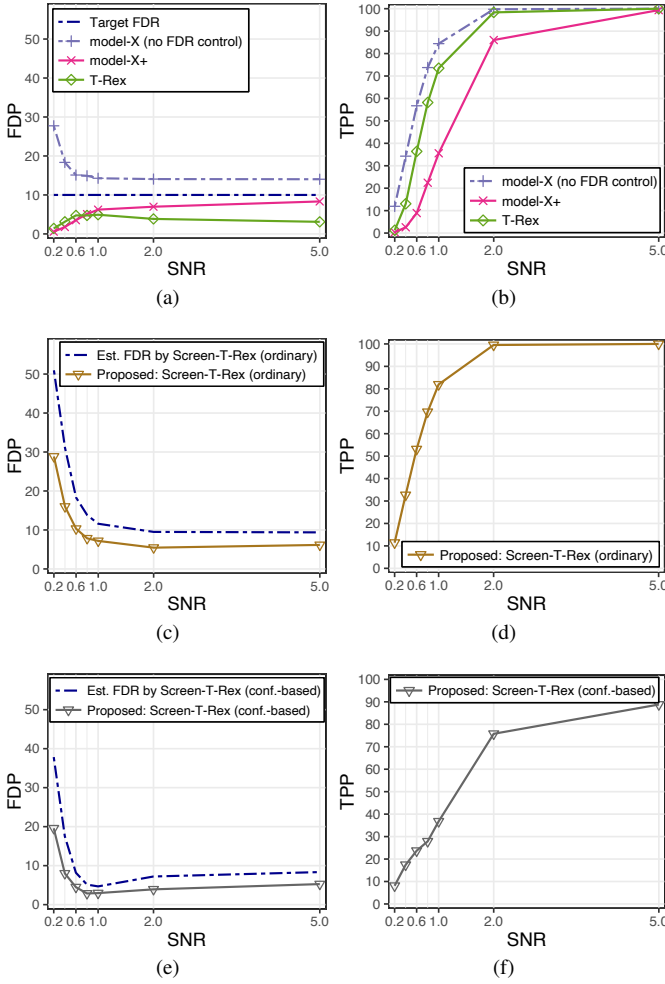


Figure 2: Figures (a) and (b) (see also [11]) display the results of the benchmark methods. We observe that the *T-Rex* selector and the *model-X* knockoff+ method control the FDR at the target level of 10%, while the *model-X* knockoff method fails to control the FDR. Figures (c) - (f) show that the proposed ordinary and the confidence-based *Screen-T-Rex* selector both control the FDR at the self-estimated levels while achieving a reasonably high TPR. The confidence-based version is capable of controlling the FDR at lower levels than the ordinary version but, in turn, achieves a lower TPR.

SNPs, i.e., SNPs that are associated with the phenotype. The results in Table I are averaged over all data sets. A discussion of the results is provided within the caption of Table I.²

VI. REAL WORLD EXAMPLE: HIV DATA

Many antiretroviral drugs are used in HIV-1 infection therapies. However, mutations may decrease the susceptibility to some drugs and, thus, lead to an increased drug resistance of the virus. Therefore, it is desired to detect mutations associated with resistance against all existing drugs to determine which drugs to use for treating HIV-1 and to develop new drugs to which mutated HIV-1 viruses are highly susceptible. In order to also compare the proposed methods

²Note that only the relative but not the absolute computation times are representative for similar settings, since, due to the energy crisis in Europe, the CPUs of the Lichtenberg High-Performance Computer of the Technische Universität Darmstadt operate at a reduced clock frequency.

Table II Results for the HIV-1 PI-type drugs: Applying Algorithm 1 with $\alpha = 3\%$, $\alpha_l = 2\%$, and $\alpha_u = 4\%$, the selected active set of the original *T-Rex* selector is only picked for ATV. For APV and IDV, the results of the confidence-based *Screen-T-Rex* selector are picked. For the remaining drugs, the results of the ordinary *Screen-T-Rex* selector are picked (compare results in Figure 3).

Drug	n	p	Target FDR	Est. FDR (ordinary)	Est. FDR (conf.-based)
Amprenavir (APV)	767	201	3 %	3.57 %	3.70 %
Atazanavir (ATV)	328	147	3 %	4.76 %	0.00 %
Indinavir (IDV)	825	206	3 %	3.33 %	3.33 %
Lopinavir (LPV)	515	184	3 %	3.85 %	0.00 %
Nelfinavir (NFV)	842	207	3 %	3.70 %	0.00 %
Ritonavir (RTV)	793	205	3 %	3.33 %	2.86 %
Saquinavir (SQV)	824	206	3 %	3.45 %	0.00 %

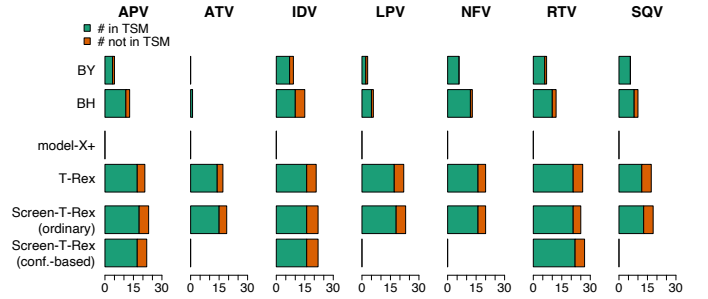


Figure 3: Number of selected mutations that are reported (green) and not reported (orange) in TSM lists for HIV-1 PI-type drugs. The *T-Rex* methods dominate the benchmark methods in terms of the number of selected mutations reported in TSM lists. Moreover, a few potentially relevant mutations that are not reported in TSM lists are detected.

against classical methods for the low-dimensional setting, we consider a low-dimensional benchmark HIV-1 data set that was described and analyzed in [23], [24] and served as a benchmark data set for the existing *fixed-X* knockoff method [6]. It can be downloaded from a Stanford University database.³ The performance of the proposed *Screen-T-Rex* selector and the benchmark methods in detecting the mutations that are associated with HIV-1 drug resistance for individual *protease inhibitor* (PI)-type drugs is assessed. The same setup as in [6] is used, i.e., the same preprocessing steps are applied and the same benchmark mutation positions from treatment-selected mutation (TSM) lists in [23] are used. Table II and Figure 3 display and discuss the results.

VII. CONCLUSION

The *Screen-T-Rex* selector, a fast FDR-controlling variable selection method for large-scale genomics biobanks, was proposed. Using the proposed method in combination with the original *T-Rex* selector, an efficient algorithm for conducting thousands of large-scale GWAS was proposed. In our future research, we will apply the proposed methods to genomics data from the UK biobank.

³URL (last access: 31st January 2023):

https://hivdb.stanford.edu/pages/published_analysis/genophenoPNAS2006/

REFERENCES

- [1] E. Uffelmann, Q. Q. Huang, N. S. Munung, J. De Vries, Y. Okada, A. R. Martin, H. C. Martin, T. Lappalainen, and D. Posthuma, "Genome-wide association studies," *Nat. Rev. Methods Primers*, vol. 1, no. 1, p. 59, 2021.
- [2] C. Sudlow, J. Gallacher, N. Allen, V. Beral, P. Burton, J. Danesh, P. Downey, P. Elliott, J. Green, M. Landray *et al.*, "UK Biobank: An open access resource for identifying the causes of a wide range of complex diseases of middle and old age," *PLOS Med.*, vol. 12, no. 3, p. e1001779, 2015.
- [3] A. Buniello, J. A. L. MacArthur, M. Cerezo, L. W. Harris, J. Hayhurst, C. Malangone, A. McMahon, J. Morales, E. Mountjoy, E. Sollis *et al.*, "The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019," *Nucleic Acids Res.*, vol. 47, no. D1, pp. D1005–D1012, 2019.
- [4] Y. Benjamini and Y. Hochberg, "Controlling the false discovery rate: a practical and powerful approach to multiple testing," *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, vol. 57, no. 1, pp. 289–300, 1995.
- [5] Y. Benjamini and D. Yekutieli, "The control of the false discovery rate in multiple testing under dependency," *Ann. Statist.*, vol. 29, no. 4, pp. 1165–1188, 2001.
- [6] R. F. Barber and E. J. Candès, "Controlling the false discovery rate via knockoffs," *Ann. Statist.*, vol. 43, no. 5, pp. 2055–2085, 2015.
- [7] E. J. Candès, Y. Fan, L. Janson, and J. Lv, "Panning for gold: 'model-X' knockoffs for high dimensional controlled variable selection," *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, vol. 80, no. 3, pp. 551–577, 2018.
- [8] R. F. Barber, E. J. Candès, and R. J. Samworth, "Robust inference with knockoffs," *Ann. Statist.*, vol. 48, no. 3, pp. 1409 – 1431, 2020.
- [9] R. F. Barber and E. J. Candès, "A knockoff filter for high-dimensional selective inference," *Ann. Statist.*, vol. 47, no. 5, pp. 2504–2537, 2019.
- [10] M. Sesia, C. Sabatti, and E. J. Candès, "Gene hunting with hidden markov model knockoffs," *Biometrika*, vol. 106, no. 1, pp. 1–18, 2019.
- [11] J. Machkour, M. Muma, and D. P. Palomar, "The terminating-random experiments selector: Fast high-dimensional variable selection with false discovery rate control," *arXiv preprint arXiv:2110.06048*, 2022. [Online]. Available: <https://arxiv.org/abs/2110.06048>
- [12] —, "False discovery rate control for grouped variable selection in high-dimensional linear models using the T-Knock filter," in *30th Eur. Signal Process. Conf. (EUSIPCO)*, 2022, pp. 892–896.
- [13] J. Machkour, S. Tien, D. P. Palomar, and M. Muma, *T-Rex Selector: T-Rex Selector: High-Dimensional Variable Selection & FDR Control*, 2022, R package version 0.0.1. [Online]. Available: <https://CRAN.R-project.org/package=T-RexSelector>
- [14] —, *tlars: The T-LARS Algorithm: Early-Terminated Forward Variable Selection*, 2022, R package version 0.0.1. [Online]. Available: <https://CRAN.R-project.org/package=tlars>
- [15] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, "Least angle regression," *Ann. Statist.*, vol. 32, no. 2, pp. 407–499, 2004.
- [16] R. Tibshirani, "Regression shrinkage and selection via the lasso," *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, vol. 58, no. 1, pp. 267–288, 1996.
- [17] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, vol. 67, no. 2, pp. 301–320, 2005.
- [18] H. Zou, "The adaptive lasso and its oracle properties," *J. Amer. Statist. Assoc.*, vol. 101, no. 476, pp. 1418–1429, 2006.
- [19] B. Efron and R. J. Tibshirani, *An introduction to the bootstrap*. CRC Press, 1994.
- [20] A. C. Davison and D. V. Hinkley, *Bootstrap methods and their application*. Cambridge University Press, 1997.
- [21] Z. Su, J. Marchini, and P. Donnelly, "HAPGEN2: simulation of multiple disease SNPs," *Bioinformatics*, vol. 27, no. 16, pp. 2304–2305, 2011.
- [22] The International HapMap 3 Consortium, "Integrating common and rare genetic variation in diverse human populations," *Nature*, vol. 467, no. 7311, pp. 52–58, 2010.
- [23] S.-Y. Rhee, W. J. Fessel, A. R. Zolopa, L. Hurley, T. Liu, J. Taylor, D. P. Nguyen, S. Slome, D. Klein, M. Horberg *et al.*, "Hiv-1 protease and reverse-transcriptase mutations: correlations with antiretroviral therapy in subtype b isolates and implications for drug-resistance surveillance," *J. Infect. Dis.*, vol. 192, no. 3, pp. 456–465, 2005.
- [24] S.-Y. Rhee, J. Taylor, G. Wadhera, A. Ben-Hur, D. L. Brutlag, and R. W. Shafer, "Genotypic predictors of human immunodeficiency virus type 1 drug resistance," *Proc. Natl. Acad. Sci. U.S.A.*, vol. 103, no. 46, pp. 17 355–17 360, 2006.