

1 Hobotnica: exploring molecular 2 signature quality

3 Alexey Stupnikov^{1,6*}, Alexey Sizykh¹, Alexander Favorov^{2,5}, Bahman Afsari², Sarah
4 Wheelan², Luigi Marchionni³, Yulia A. Medvedeva^{1,4,6*}

*For correspondence:

aleksej.stupnikov@phystech.edu
(FMS); ju.medvedeva@gmail.com
(FS)

5 ¹Moscow Institute of Physics and Technology, Moscow, Russia; ²Johns Hopkins
6 University, Baltimore, USA; ³Weill Cornell Medicine, New York, USA; ⁴Center of
7 Biotechnology RAS, Moscow, Russia; ⁵Vavilov Institute for General Genetics RAS,
8 Moscow, Russia; ⁶National Medical Research Center for Endocrinology, Moscow, Russia

10 **Abstract** A Molecular Features Set (MFS), is a result of vast diversity of bioinformatics pipelines.
11 In case when MFS is used for further analysis to distinguish between phenotypes, it is often
12 referred to as a signature. Lack of the “gold standard” for most experimental data modalities
13 makes it hard to provide valid estimation for a particular MFS’s quality. Yet, this goal can partially
14 be achieved by analyzing inner-sample Distance Matrix (DM) and their power to distinguish
15 between phenotypes.
16 The quality of a DM can be assessed by summarizing its power to quantify the differences of
17 inner-phenotype and outer-phenotype distances. This estimation of the DM quality can be
18 construed as a measure of the MFS’s quality.
19 Here we propose **Hobotnica**, an approach to estimate MFS’s quality by their ability to stratify
20 data, and assign them significance scores, that allows for collating various signatures and
21 comparing their quality for contrasting groups.

23 Introduction

24 A signature based on a predefined Molecular Features Set (MFS), which is designed to distinguish
25 biological conditions or phenotypes from each other — is one of major concepts of bioinformatics
26 and precision medicine. In this context, signatures typically originate from MFS from contrasting
27 experimental data of two or more sample groups, which differ phenotypically. These MFS incor-
28 porate information on the differences between the groups. The nature of the MFS depends on
29 the modality of the original data. For instance, the MFS provided by the Differential Gene Express-
30 sion approach is a list of Differentially Expressed Genes (DEG); Differential Methylation analysis
31 provides Differentially Methylated Cytosines or Regions (DMC and DMR) as MFS, etc.

32 A significant number of mutational, expression and methylation-based signatures have recently
33 been published and they are actively used in Research and Transnational Medicine. Examples
34 of expression-based signatures involve genesets for clinical prognosis (e.g. PAM50 (*Parker et al.*
35 *(2009)*), MammaPrint (*Cardoso et al. (2016)*) for Breast Cancer), for pathways and gene enrichment
36 analysis (e.g. MsigDB collections (*Subramanian et al. (2005)*)), for drug re-purposing (e.g. LINCS
37 project(*Liu et al. (2015)*)).

38 Direct quality assessment for MFS is currently hardly possible, since there are no ‘gold standard’
39 datasets where active Molecular Features are explicitly known. In this manuscript, we propose a
40 novel approach - **Hobotnica** - that allows for measurement of MFS quality by addressing the key
41 property of the signature, namely, its quality for data stratification.

42 Hobotnica leverages the quality of Distance Matrices obtained from any source in order to as-

43 assess quality of the MFS from any data modality compared to a random MFS. In this study, we
44 demonstrate its application on transcriptomic signatures.

45 Results

46 Approach

47 The Hobotnica approach is as follows: For a given data set W and a given Molecular Features Set
48 (S) we derive the inter-sample distance matrix ($DM(S, W)$). Then we assess the quality of DM
49 (and, thus, of S) with a summarizing function ($\alpha(DM(S)) = \alpha(DM(S), Y)$ or by abuse of notation
50 $\alpha(DM(S))$) where (Y) represents the labels of samples.

51 We desire the function α to gauge if the inner-class samples are closer to each other than to
52 outer-class samples. If no difference exists from one class to another, α must be close to zero and
53 as the difference grows, α grows. In ideal case of a perfect separation, α reaches its maximum at
54 1:

- 55 • $\alpha \in [0, 1]$
- 56 • $\alpha \rightarrow 1 \Leftrightarrow$ High groups stratification quality
- 57 • $\alpha \rightarrow 0 \Leftrightarrow$ Low groups stratification quality

58 Under the Null hypothesis of Hobotnica ((H_0)), no significant difference exists between $\alpha(S)$ and
59 the α of an equal-sized general random set. On the contrary, the Alternative ((H_A)) hypothesizes that
60 S generates higher α than most random S' of the same size. To estimate a Null distribution for
61 *Hobotnica's* α , we applied a permutation test. As our default options, we use Kendall distance as
62 the distance measure and Mann-Whitney-Wilcoxon test as the summarizing function.

63 Validation

64 To validate our approach in the first case study we extracted RNA-seq expression dataset for
65 Prostate Cancer from TCGA on counts level (*Rahman et al. (2015)*). As Molecular Feature Sets we
66 recruited C2 collection of molecular signatures from MSigDB (*Subramanian et al. (2005)*, *Liberzon*
67 *et al. (2011)*) that contains a number of Prostate-related genesets. For the second case study we
68 took PAM50 molecular signature, designed for various Breast Cancer types classification, and ap-
69 plied it to several datasets (*Marusyk et al. (2016)*)(*Daemen et al. (2013)*)(*Costello et al. (2014)*)(*Rah-*
70 *man et al. (2015)*)(*Varley et al. (2014)*). In both cases, the counts were normalised to cpm. For each
71 geneset H-score and its p-value with BH correction were computed.

72 Prostate-related C2 genesets clearly demonstrated highest values of H-score and sufficient sta-
73 tistical significance (Fig.1.A), as well as data stratification (Fig.1.B), which is expected for Prostate
74 Cancer vs Control contrast. Genesets not attributed to Prostate Cancer related processes did not
75 achieve statistical significant p-values. (Table1).

76 PAM50 signature evidently separates samples in for GSE48216 dataset (Fig.1.C). H-scores for
77 random genesets for the same dataset are significantly lower than an H-score for PAM50 (Fig.1.D).
78 Clearly, PAM50 signature demonstrates high quality of stratification for the samples of various
79 Breast Cancer datasets with high H-score values and statistically significant p-values (Table 2).

80 Thus, in the first case study, Prostate Cancer related genesets from C2 collection, when applied
81 to Prostate Cancer dataset, delivered highest H-scores and most significant and p-values proved to
82 demonstrate best scores and performance. Likewise, in the second case study, PAM50 expression
83 signature applied to several heterogeneous Breast Cancer datasets delivered high H-score values
84 along with significant scores of p-values.

85 Application

86 An important question that researches often face is establishing the optimal size of the retrieved
87 signature. The exact number of genes to be retrieved from the set of all significant genes is an
88 important parameter that is essential for signature's application. To explore the optimal size of DE

Table 1. 10 C2-CGP Gene Signatures with highest H-scores

Signature	H-score	p-value
TOMLINS_PROSTATE_CANCER	0.795	0.025
WALLACE_PROSTATE_CANCER	0.747	0.025
OUYANG_PROSTATE_CANCER_PROGRESSION	0.745	0.025
LIU_PROSTATE_CANCER	0.735	0.025
PIEPOLI_LGI1_TARGETS	0.724	0.059
SMID_BREAST_CANCER_RELAPSE_IN_LIVER	0.712	0.164
TIMOFEEVA_GROWTH_STRESS_VIA_STAT1	0.708	0.240
GENTILE_UV_LOW_DOSE	0.705	0.308
JOHANSSON_BRAIN_CANCER_EARLY_VS_LATE	0.701	0.377
HOWLIN_CITED1_TARGETS_1	0.700	0.377

Table 2. PAM50 results

GEO Accession	Sample size	Groups in dataset	H-score	p-value
GSE58135	168	6	0.772	7e-4
GSE62944	1067	5	0.8892	0.0003
GSE48216	46	3	0.8567	0.0003
GSE80333	10	3	0.9765	0.0003

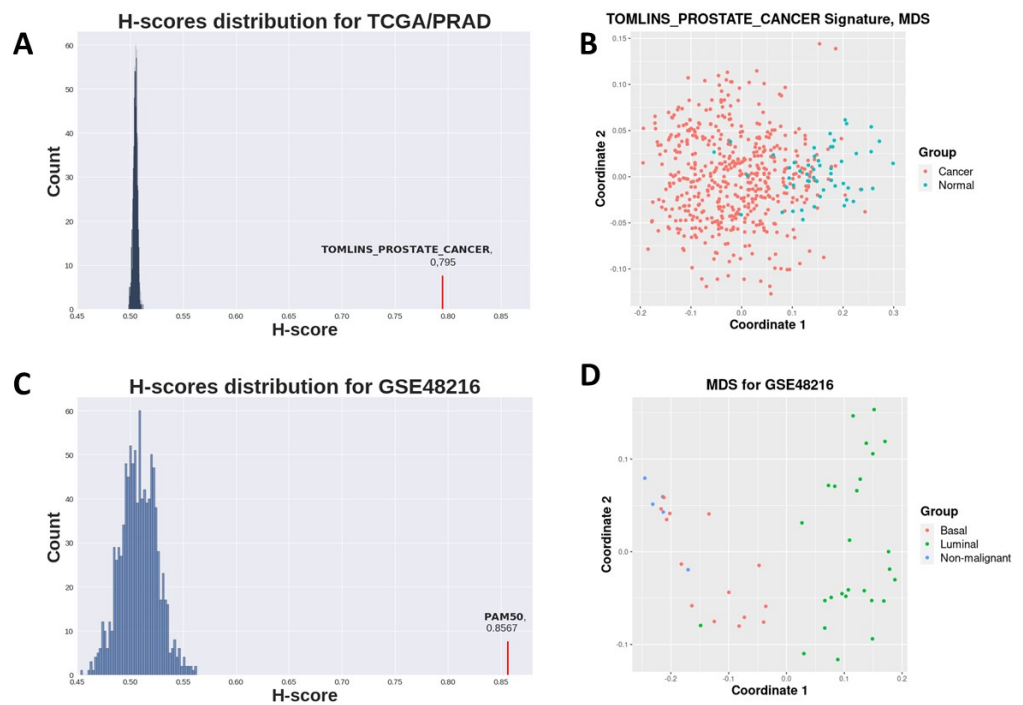


Figure 1. A : Distribution of H-scores for random genesets (blue) on TCGA Prostate Cancer vs Normal dataset (see Tab.1) and Tomlins prostate geneset H-score (red). B: MDS for TCGA Prostate demonstrates samples separation with Tomlins geneset. C: Distribution of H-scores for random genesets (blue) on GSE48216 Breast Cancer dataset (see Tab.2) and PAM50 geneset H-score (red). D: MDS for GSE48216 Breast Cancer dataset samples separation with PAM50 geneset.

89 signature we performed Hobotnica analysis for top DE p-value ordered gene signatures of various
 90 lengths. For the reference we performed DGE analysis for Breast Cancer vs Control TCGA dataset

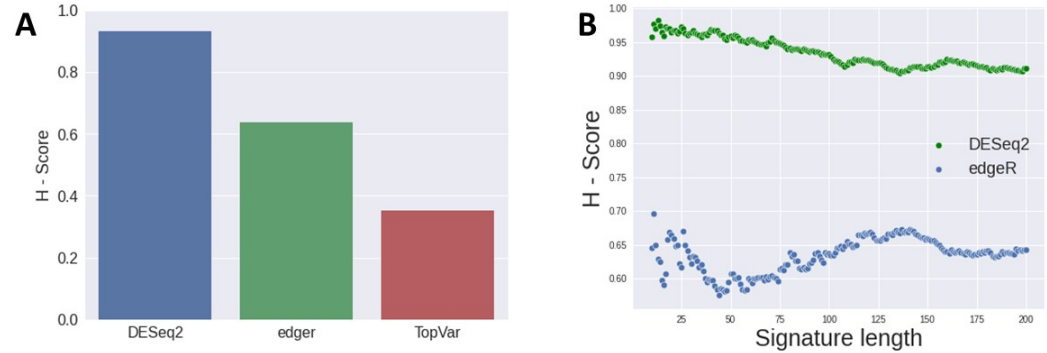


Figure 2. A: H-scores delivered by top 100 gene signatures from various DE models applied to TCGA Breast Cancer data. B: Change of H-score with the length of gene signature derived from DESeq2 and edgeR models

91 (*Rahman et al. (2015)*) with DESeq2 (*Love et al. (2014)*) and edgeR (*McCarthy et al. (2012)*). Top
 92 100 genes for each method were retrieved, as well as genes with highest variance in expression.
 93 H-scores for every signature then were computed (Fig.2.A.). For this dataset DESeq2 provided a
 94 signature with the highest quality score. Then, we calculated H-score for signatures of various
 95 lengths (Fig.2.B.) Surprisingly, the signature quality is non-monotonously dependant on the signa-
 96 ture length, i.e. number of genes in the signature. The pattern also varies for DE models. Addi-
 97 tionally, for the best-performing model, DESeq2, the a signature quality is generally declining with
 98 the length. Thus, increasing number of genes in a signature may not improve its quality, and an
 99 optimal gene signature length for the DE analysis result may be established: DESeq2 signature
 100 reaches its maximum H-score at 13 genes and edgeR at 11 in this case.

101 Discussion

102 *Hobotnica* is designed to quantitatively evaluate Molecular Feature Set's quality by their ability for
 103 data stratification from their inter-sample distance matrices, and to assess the statistical signifi-
 104 cance of the results. We demonstrated that *Hobotnica* can efficiently estimate the quality of a
 105 Molecular Signature in the context of Expression data.

106 Suggested method can be used to evaluate Molecular Feature sets of various nature: retrieved
 107 in DGE, Differential Methylation analysis, Mutation/SNV calling or Pathways analysis, as well as
 108 data modalities from other types of Differential Problem. In addition, assessing H-score values for
 109 various lengths of the same set or signature will help with its structure optimization, which may be
 110 especially important in clinical applications.

111 *Hobotnica* is available as an R package at <https://github.com/lab-medvedeva/Hobotnica-main>

112 Methods and Materials

113 Problem formalization

114 If a Molecular Feature Set (S), that presumably incorporates information on the contrast between
 115 groups of samples with known samples annotation Y in Data D is in place ($H : S$), we can com-
 116 pute Distance Matrix between samples $DM (f(S|D) \rightarrow DM)$ and then introduce a measure α
 117 ($g(DM|Y) \rightarrow \alpha$) of signature quality for Data D stratification.

$$\begin{aligned}
 &H : S \\
 &f(S|D) \rightarrow DM \\
 &g(DM|Y) \rightarrow \alpha
 \end{aligned} \tag{1}$$

118 When instead of a single GS a set of hypotheses $\{H_1 : GS_1, H_2 : GS_2, \dots, H_n : GS_n\}$ is in place,
119 for each Gene Signature GS_i corresponding Distance Matrix DM_i can be generated, and then, in
120 turn, particular value of the measure α_i :

$$\begin{cases} H_1 : S_1 \\ H_2 : S_2 \\ \dots \\ H_n : S_n \end{cases} \rightarrow \begin{cases} f(S_1|D) \rightarrow DM_1 \\ f(S_2|D) \rightarrow DM_2 \\ \dots \\ (S_n|D) \rightarrow DM_n \end{cases} \rightarrow \begin{cases} g(DM_1|A) \rightarrow \alpha_1 \\ g(DM_2|A) \rightarrow \alpha_2 \\ \dots \\ g(DM_n|A) \rightarrow \alpha_n \end{cases} . \quad (2)$$

121 Thus, for every MFS S_i from set of hypotheses $\{H_1 : S_1, H_2 : S_2, \dots, H_n : S_n\}$ H-score α_i may be
122 computed, resulting in a set $\langle \alpha_1, \alpha_2, \dots, \alpha_n \rangle$. Comparing α values allows for corresponding Feature
123 Sets qualities ranking and selecting the most informative Signatures for the Data D .

124 To assess statistical significance of each obtained H-score α_i we compute empirical *p-value* via
125 generating a distribution of H-scores for set of random MFS.

126 Availability

127 We implemented Hobotnica as an R package available at [https://github.com/lab-medvedeva/Hobotnica-](https://github.com/lab-medvedeva/Hobotnica-main)
128 [main](https://github.com/lab-medvedeva/Hobotnica-main)<https://github.com/lab-medvedeva/Hobotnica-main>. It contains an implementation of the
129 Hobotnica measure, statistical analysis for significance, and several auxiliary functions for visu-
130 alizing results and parallel processing.

131 Acknowledgements

132 We thank Frank Emmert-Streib, Leslie Cope and Elana Fertig for fruitful discussions. The study was
133 supported by Ministry of Science and Higher Education of the Russian Federation (agreement no.
134 075-15-2020-899) and by the NIH grants R01DE027809 and P30CA006973.

135 References

- 136 **Cardoso F**, van't Veer LJ, Bogaerts J, Slaets L, Viale G, Delaloge S, Pierga JY, Brain E, Causeret S, DeLorenzi M,
137 et al. 70-gene signature as an aid to treatment decisions in early-stage breast cancer. *New England Journal*
138 *of Medicine*. 2016; 375(8):717–729.
- 139 **Costello JC**, Heiser LM, Georgii E, Gönen M, Menden MP, Wang NJ, Bansal M, Hintsanen P, Khan SA, Mpindi JP,
140 et al. A community effort to assess and improve drug sensitivity prediction algorithms. *Nature biotechnology*.
141 2014; 32(12):1202–1212.
- 142 **Daemen A**, Griffith OL, Heiser LM, Wang NJ, Enache OM, Sanborn Z, Pepin F, Durinck S, Korkola JE, Griffith M,
143 et al. Modeling precision treatment of breast cancer. *Genome biology*. 2013; 14(10):1–14.
- 144 **Liberzon A**, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular signatures
145 database (MSigDB) 3.0. *Bioinformatics*. 2011 05; 27(12):1739–1740. [https://doi.org/10.1093/bioinformatics/](https://doi.org/10.1093/bioinformatics/btr260)
146 [btr260](https://doi.org/10.1093/bioinformatics/btr260), doi: 10.1093/bioinformatics/btr260.
- 147 **Liu C**, Su J, Yang F, Wei K, Ma J, Zhou X. Compound signature detection on LINCS L1000 big data. *Molecular*
148 *BioSystems*. 2015; 11(3):714–722.
- 149 **Love MI**, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with
150 DESeq2. *Genome biology*. 2014; 15(12):1–21.
- 151 **Marusyk A**, Tabassum DP, Janiszewska M, Place AE, Trinh A, Rozhok AI, Pyne S, Guerriero JL, Shu S, Ekram M,
152 et al. Spatial proximity to fibroblasts impacts molecular features and therapeutic sensitivity of breast cancer
153 cells influencing clinical outcomes. *Cancer research*. 2016; 76(22):6495–6506.
- 154 **McCarthy DJ**, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq experiments with
155 respect to biological variation. *Nucleic Acids Research*. 2012 01; 40(10):4288–4297. doi: 10.1093/nar/gks042.
- 156 **Parker JS**, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, Davies S, Fauron C, He X, Hu Z, et al. Supervised
157 risk predictor of breast cancer based on intrinsic subtypes. *Journal of clinical oncology*. 2009; 27(8):1160.

- 158 **Rahman M**, Jackson LK, Johnson WE, Li DY, Bild AH, Piccolo SR. Alternative preprocessing of RNA-Sequencing
159 data in The Cancer Genome Atlas leads to improved analysis results. *Bioinformatics*. 2015; 31(22):3666–3672.
- 160 **Subramanian A**, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR,
161 Lander ES, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide
162 expression profiles. *Proceedings of the National Academy of Sciences*. 2005; 102(43):15545–15550.
- 163 **Varley KE**, Gertz J, Roberts BS, Davis NS, Bowling KM, Kirby MK, Nesmith AS, Oliver PG, Grizzle WE, Forero A,
164 et al. Recurrent read-through fusion transcripts in breast cancer. *Breast cancer research and treatment*.
165 2014; 146(2):287–297.