

1 **treeClust improves protein co-regulation analysis due to robust selectivity for**  
2 **close linear relationships**

3 Georg Kustatscher<sup>1</sup>, Piotr Grabowski<sup>2</sup>, Juri Rappsilber<sup>1,2#</sup>

4 <sup>1</sup> Wellcome Centre for Cell Biology, University of Edinburgh, Edinburgh EH9 3BF, UK

5 <sup>2</sup> Bioanalytics, Institute of Biotechnology, Technische Universität Berlin, 13355 Berlin, Germany

6 # Communicating author: [juri.rappsilber@ed.ac.uk](mailto:juri.rappsilber@ed.ac.uk)

7 **Gene co-expression analysis is a widespread method to identify the potential biological**  
8 **function of uncharacterised genes. Recent evidence suggests that proteome profiling**  
9 **may provide more accurate results than transcriptome profiling. However, it is unclear**  
10 **which statistical measure is best suited to detect proteins that are co-regulated. We have**  
11 **previously shown that expression similarities calculated using treeClust, an**  
12 **unsupervised machine-learning algorithm, outperformed correlation-based analysis of a**  
13 **large proteomics dataset. The reason for this improvement is unknown. Here we**  
14 **systematically explore the characteristics of treeClust similarities. Leveraging synthetic**  
15 **data, we find that tree-based similarities are exceptionally robust against outliers and**  
16 **detect only close-fitting, linear protein - protein associations. We then use proteomics**  
17 **data to demonstrate that both of these features contribute to the improved performance**  
18 **of treeClust relative to Pearson, Spearman and robust correlation. Our results suggest**  
19 **that, for large proteomics datasets, unsupervised machine-learning algorithms such as**  
20 **treeClust may significantly improve the detection of biologically relevant protein - protein**  
21 **associations relative to correlation metrics.**

22 **INTRODUCTION**

23 Genes with related biological functions tend to be active in the same biological conditions. This  
24 is the basis of gene coexpression analysis, a method that predicts the function of unknown  
25 genes by comparing their expression profiles to those of well-studied genes (1–5). A typical  
26 coexpression study detects gene activity by measuring mRNA abundances of many genes in a  
27 range of biological samples or conditions. In a second step, the similarity of expression profiles,  
28 i.e. the extent of coexpression between any two genes, is determined by correlation analysis.  
29 Finally, pairwise coexpression coefficients are aggregated into a gene coexpression network,  
30 through which any uncharacterised genes in the dataset may become associated with clusters  
31 of genes of well-defined biological functions.

32 Many variations of this basic approach have been developed over the past two decades  
33 (6). For example, several coexpression measures have been explored as alternatives to  
34 Pearson's correlation, including Spearman's correlation, Biweight midcorrelation, Mutual  
35 Information and simple regression models (7, 8). No single measure appears to be superior for

36 every dataset, as the optimal choice of measure depends on various characteristics of a given  
37 dataset, such as the frequency of outliers and missing values.

38 A recent, fundamental change to the expression profiling setup was made possible by  
39 improvements in the field of quantitative proteomics: the use of protein abundances rather than  
40 mRNAs as readout for gene activity. This increases the accuracy of gene function prediction,  
41 because protein abundances are better indicators of gene function than mRNA levels, at least in  
42 human (9–11) and mouse (12). We have recently reported ProteomeHD, a dataset that  
43 quantifies the response of 10,323 human proteins to 294 biological perturbations using  
44 isotope-labelling mass spectrometry (13) (<https://www.biorxiv.org/content/10.1101/582247v1>, in  
45 revision). ProteomeHD is a heterogeneous dataset, incorporating a wide range of perturbation  
46 experiments from different laboratories, such as inhibitor treatments, differentiation time courses  
47 and cancer cell line comparisons. We compared different coexpression measures for their ability  
48 to detect proteins that are co-regulated in response to these perturbations. Surprisingly, we  
49 found that the unsupervised machine-learning algorithm treeClust (14, 15) provided a striking  
50 improvement over established correlation-based metrics. However, the reason for this  
51 improvement remained unclear, because treeClust is a novel algorithm that works in a  
52 fundamentally different way to previously used coexpression metrics.

53 The treeClust algorithm uses recursive partitioning (16–18) to create decision trees.  
54 Such trees are normally used for supervised classification or regression tasks. In contrast,  
55 treeClust uses decision trees to calculate a dissimilarity measure in an unsupervised manner. To  
56 do so, treeClust first creates a number of decision trees aimed to dissect the dataset by growing  
57 one decision tree for each variable in turn, using it as response variable and all remaining  
58 variables as predictors. This part of the algorithm is reminiscent of an established approach to  
59 impute missing data which uses Random Forests (19). However, in a second step, treeClust  
60 calculates dissimilarities between any two observations based on the proportion of trees in  
61 which they land in different leaves, using Gower's distance (20). While treeClust dissimilarities  
62 appear to perform well in practical applications (21, 22), some of their basic properties remain  
63 unclear, especially in the context of gene expression analysis. For example, do treeClust  
64 dissimilarities capture linear or non-linear associations? How is treeClust performance affected  
65 by missing data, outlier data and noise? Here, we set out to systematically address these  
66 questions and benchmark treeClust performance on both synthetic and real proteomics data.

## 67 **EXPERIMENTAL PROCEDURES**

### 68 **General data analysis and availability**

69 All data processing and analysis has been performed using R version 3.5.1 (23). All data and R  
70 scripts required to reproduce the results of this manuscript are available in the following GitHub  
71 repository: <https://github.com/Rappsilber-Laboratory/treeClust-benchmarking>.

72 The R package data.table (24) was used for fast data processing. Figures were prepared  
73 using ggplot2 (25), gridExtra (26), cowplot (27) and viridis (28).

## 74 **Generation of synthetic datasets**

75 Synthetic datasets were generated using a custom function in R. The function populates a table  
76 with values that are randomly sampled from a normal distribution, but includes a user-specified  
77 number of observations that have a defined linear relationship with each other. The following  
78 properties of the thus created datasets can be manipulated: number of variables (i.e. samples or  
79 experiments), number of observations (i.e. proteins), percentage of protein pairs that should  
80 have a linear relationship, percentage of outlier data, percentage of missing values and the  
81 extent of scatter around the regression line (i.e. biological or measurement noise). Outlier data  
82 points are created by random sampling from a broader normal distribution than the rest of the  
83 data.

84 In addition to positive linear relationships ( $y \sim x$ ), we tested relationships that were  
85 exponential ( $y \sim e^x$ ), logistic ( $y \sim 4 / (1 + e^{-5x})$ ) and quadratic ( $y \sim x^2$ ), as well as linearly  
86 anti-correlated ( $y \sim -x$ ).

## 87 **Real biological datasets**

88 ProteomeHD has been documented in detail elsewhere (13). In short, it is a data matrix  
89 consisting of 10,323 proteins whose abundance changes in response to 294 biological  
90 perturbations have been determined by quantitative mass spectrometry, using stable isotope  
91 labelling by amino acids in cell culture (29). To distinguish between genuine, biologically relevant  
92 protein - protein associations (true positives) and likely false positive interactions we used a gold  
93 standard based on the Reactome database (30), which was also described previously (13).

## 94 **Comparison of coexpression measures**

95 Pearson's correlation coefficients and Spearman's rank correlation coefficients were calculated  
96 using R base functions (23). Biweight Midcorrelation (bicor) was calculated with default settings  
97 using the R package WGCNA (31, 32). TreeClust dissimilarities were calculated using the R  
98 package treeClust (14, 15), with the d.num parameter set to 2. When applying treeClust to  
99 ProteomeHD rather than synthetic data, we set the rpart complexity parameter to 0.105 and the  
100 treeClust serule parameter to 1.8. These settings had been optimised previously for  
101 ProteomeHD (13), providing approximately a 10% performance improvement over default  
102 values when assessed against the Reactome gold standard.

103 Performance of coexpression measures was compared by precision - recall (PR)  
104 analysis using the R package PRROC (33). True positive (linear or nonlinear) and false positive  
105 (random) associations for the PR analyses were known a priori for synthetic data and annotated  
106 using the Reactome gold standard for ProteomeHD. To test the impact of various data  
107 characteristics, synthetic dataset were generated in triplicate and the result is shown as the  
108 average area under the PR curves, with error bars indicating the standard error of the mean. No  
109 replicates were used for the combinatorial testing of two dataset characteristics (Fig. 2C, G and  
110 H).

## 111 **Model fitting in real proteomics data**

112 Base R functions were used to fit and analyse linear models for pairs of proteins in  
113 ProteomeHD. Fold-changes of each protein pair were rescaled to fall between 0 and 1 before  
114 fitting the model. Outliers were defined as data points with absolute studentized residuals or a  
115 Mahalanobis distance larger than 2. Non-linear models were fit using nonlinear least squares.  
116 Exponential models ( $y \sim a + \exp(b)^x$ ) and logistic models ( $y \sim a / (1 + e^{-b(x-c)})$ ) were said to  
117 outperform the corresponding linear model ( $y \sim a + bx$ ) if their residual sum of squares (RSS)  
118 was at least 10% smaller.

## 119 **RESULTS AND DISCUSSION**

### 120 **A coexpression take on Anscombe's quartet**

121 In the protein co-regulation analysis of ProteomeHD treeClust outperformed common  
122 coexpression measures: Pearson's correlation coefficient (PCC), Spearman's rank correlation  
123 ( $\rho$ ) and Biweight midcorrelation (bicor). To explore possible reasons for this we used  
124 Anscombe's quartet (34) as a starting point. These four 11-point datasets illustrate several key  
125 issues that can negatively affect the performance of Pearson's correlation (Fig. 1A). For  
126 example, PCC can falsely identify a linear correlation when there is a non-linear relationship  
127 between two variables. In addition, outlier data located far off the regression line can lead to an  
128 underestimation of the correlation. Similarly, outliers can cause a high PCC when in fact no  
129 correlation between two variables exists at all. Spearman's  $\rho$  and bicor are also affected by  
130 these issues, albeit to a much lesser extent than PCC (Fig. 1A).

131 We then asked how treeClust deals with Anscombe's quartet. However, it is not possible  
132 to simply calculate treeClust dissimilarities for Anscombe's four variable pairs. This is because,  
133 being a machine-learning algorithm, treeClust requires an input dataset with many variables in  
134 order to build informative decision trees. Therefore, we created a series of synthetic datasets  
135 that allow us to systematically assess the properties of treeClust dissimilarities and compare  
136 them to the properties of common coexpression measures. For example, we created a synthetic  
137 dataset consisting of 100 variables (experiments, samples or biological conditions) and 200  
138 observations (proteins). The dataset is built in such a way that 99.5% of the resulting pairwise  
139 "protein - protein" associations are random, i.e. values for both proteins are random samples of  
140 a normal distribution (Fig. 1B). The remaining 0.5% pairs are designed to have a clearly defined,  
141 linear relationship across the 100 "experiments". These pairs have a PCC close to 1, which  
142 clearly sets them apart from the distribution of the random pairs (Fig. 1B). On such a synthetic  
143 dataset, treeClust generates 100 decision trees that assign very different dissimilarities to  
144 random and linear associations (Fig. 1B). Indeed, in this simple best-case scenario, all of the  
145 four tested coexpression measures separate random and defined pairs perfectly, resulting in  
146 precision - recall curves with an area of 1 (Fig. 1C).

147 Note that in this manuscript we show treeClust similarities rather than dissimilarities  
148 (similarity = 1 - dissimilarity), in order to make the comparison with correlation metrics more  
149 intuitive (Fig. 1B).

### 150 **Linear vs non-linear relationships**

151 We then proceeded to modify various properties of the synthetic datasets and assessed how  
152 they affect treeClust. First, we asked which types of associations are detected by treeClust. For  
153 example, we replaced the 0.5% linearly correlated pairs with exponential relationships (Fig. 1D).  
154 This does not affect Pearson, Spearman or robust correlation, which still yield an area under the  
155 precision - recall curve (AUPRC) of 1. Although exponentially related pairs receive lower  
156 correlation coefficients than linear ones, their coefficients are still much higher than those of  
157 random pairs. Surprisingly and in stark contrast to the correlation measures, treeClust does not  
158 detect exponential relationships at all, yielding an AUPRC of 0 (Fig. 1D). We obtained the same  
159 result for logistic relationships (Fig. 1D). None of the coexpression measures detects quadratic  
160 relationships. Finally, we tested if treeClust detects negative linear associations, i.e.  
161 anti-correlation. We find that treeClust only partially separates anti-correlated from random  
162 associations, suggesting that low treeClust similarities indicate a lack of correlation rather than  
163 anti-correlation (Fig. 1E).

164 We conclude that, in the conditions tested here, treeClust specifically captures positive  
165 linear “protein - protein” associations. This property could be explained by the fact that treeClust  
166 dissimilarities reflect how often two observations land in the same decision tree leaf. A split in a  
167 decision tree is less likely to separate two linearly associated proteins than exponentially related  
168 or anti-correlated proteins. For example, if protein X1 is upregulated 1.5-fold in a given  
169 experiment, a linearly related protein X2 may be upregulated 1.6-fold. These proteins would only  
170 land in different leaves if a split occurred between 1.5 and 1.6. However, if protein X2 was  
171 exponentially related to X1 it may be upregulated by 4.6-fold, increasing the margin within which  
172 a split could occur such that the two proteins land in different leaves. Similarly, if two proteins  
173 are anticorrelated they rarely end up in the same leaf and thus cannot be flagged up by  
174 treeClust as being linked.

### 175 **Size and structure of the dataset**

176 We next investigated how basic data characteristics such as the number of variables and  
177 observations affect treeClust. In principle, treeClust performance is expected to improve with the  
178 amount of data it is presented with, because more data may allow treeClust to build more  
179 informative decision trees and thus learn better to distinguish between random and genuine  
180 linear associations. To test this, we constructed a series of synthetic datasets with increasing  
181 dimensions.

182 First, we increase the number of variables, i.e. samples (Fig. 2A). Under our test  
183 conditions treeClust requires around 40 samples to reach optimal performance, i.e. an AUPRC  
184 of 1. In contrast, the three correlation metrics only need ~15 samples to reliably identify all

185 genuinely correlated proteins (Fig. 2A). A likely explanation for this difference is that treeClust  
186 builds one decision tree per sample, so increasing the number of variables also increases the  
187 number of decision trees and thus the reliability of the resulting dissimilarities.

188 Next, we modify the percentage of linear associations (i.e. coexpressed proteins) in the  
189 dataset (Fig. 2B). This has no impact on the correlation measures but affects treeClust  
190 performance, suggesting that the more genuine associations are present in the data the easier it  
191 is for treeClust to learn to identify them. Notably, these two data size characteristics are  
192 interdependent: increasing the number of samples compensates for a smaller percentage of  
193 defined associations and vice versa (Fig. 2C).

194 Third, we test the impact of having more observations (proteins) in a dataset. We find  
195 that increasing the number of proteins to  $> 1,000$  is sufficient for optimal treeClust performance  
196 even if only 20 samples are available (Fig. 2D). Larger and more complex input data generally  
197 improve the performance of machine-learning algorithms. While increasing observations will not  
198 affect the number of decision trees, the increased complexity of the input data allows treeClust  
199 to create more informative trees.

200 In summary, these results show that optimal treeClust performance requires a dataset of  
201 a certain size and structure, for example  $\geq 50$  variables,  $\geq 1000$  proteins and  $\geq 0.4\%$  linear  
202 associations. For smaller datasets, for example in the range of 20 variables and 500  
203 observations, traditional correlation-based measures may be better suited for coexpression  
204 analysis.

## 205 **Missing values**

206 Proteomics data often contain a large number of missing values. In this case, correlation metrics  
207 simply focus on those variables that have been measured for both proteins. The decision trees  
208 used by treeClust handle missing values through surrogate splits (14, 17), an approach that is  
209 generally considered to be sensible only if missing values are sparse. To evaluate the impact of  
210 missing values on treeClust performance we randomly introduce missing values in synthetic  
211 data. In a dataset with 50 variables and 500 observations, introducing 10% - 15% of missing  
212 values has no ill-effect on treeClust performance (Fig. 2E). Beyond that, missing values quickly  
213 become detrimental for treeClust. In contrast, they do not pose a problem for correlation metrics  
214 as long as a sufficient number of common pairwise measurements remain available (Fig. 2E).  
215 However, we find that the impact of missing values depends on the overall dimensions of the  
216 input data. For example, with a dataset of 100 samples and 1000 proteins treeClust can already  
217 tolerate 20% missing values (Fig. 2F). Consequently, we systematically explore the impact of  
218 missing values depending on the number of samples or proteins. We observe that for large  
219 datasets treeClust performance does not decrease even if 40% of all values are missing (Fig.  
220 2G, H).

## 221 **Goodness-of-fit**

222 We then asked how “tight” coexpression of two hypothetical proteins needs to be for treeClust to  
223 detect it. To this end, we increased the dispersion / scatter of values around the linear

224 associations (Fig. 2I). Within the range of parameters tested, increasing dispersion had no  
225 appreciable effect on the performance of the three correlation metrics (Fig. 2J). In contrast,  
226 treeClust detects only very close, well-fitting linear associations. As for non-linear relationships,  
227 the explanation for this behaviour may lie in the likelihood of decision tree splits occurring  
228 between two observations. A larger scatter around the regression lines signifies a larger  
229 difference between two proteins and therefore an increased probability for them to land in  
230 different leaves.

### 231 **Outlier data points**

232 Finally, we assessed the impact of outlier data on treeClust dissimilarities. Introducing outlier  
233 data points in a synthetic dataset confirms the well known error-proneness of Pearson's  
234 correlation in the presence of outliers (Fig. 2K, L). As expected from Anscombe's examples, the  
235 AUPRC for Pearson's correlation is halved if around 25% of the measurements are outliers.  
236 Spearman and Biweight midcorrelation, which are less susceptible to outliers, handle this level  
237 of outliers in our test set without performance decrease (Fig. 2K, L). However, treeClust is  
238 exceptionally robust against outlier measurements, even in comparison to Spearman's rho and  
239 bicor. In the synthetic dataset, treeClust performance is completely unaffected by up to 75%  
240 outliers. Therefore, treeClust can detect an association between two synthetic proteins if only  
241 25% of the actual measurements show a strong linear relationship.

### 242 **Applying the lessons from synthetic data to real proteomics experiments**

243 The synthetic datasets revealed several marked differences between treeClust and traditional  
244 correlation-based coexpression measures. One potential disadvantage of treeClust  
245 dissimilarities is that they can only be calculated accurately for datasets that fulfill certain  
246 requirements on size and structure, including the number of experiments, proteins, percentage  
247 of coexpressed protein pairs and missing values.

248 A dataset like ProteomeHD is well within the margins of optimal treeClust performance  
249 identified by the synthetic data. We applied treeClust to 5,013 proteins in ProteomeHD that had  
250 been observed in at least a third of the 294 samples. This subset of ProteomeHD contains 35%  
251 missing values and a sufficient percentage of genuine linear protein-protein associations. The  
252 latter is estimated based on the observation that 3% of all protein pairs in this dataset have  
253 strong and significant Pearson's correlations ( $PCC > 0.5$ , Bonferroni adjusted  $p$ -values  $< 1e-6$ ).  
254 This is well in excess of the 0.5% margin determined on synthetic data, even after accounting  
255 for a reasonable fraction of potential false-positives.

256 Using the synthetic data we identified two potential reasons for the improved  
257 co-regulation analysis of ProteomeHD. First, treeClust detects exclusively close, linear  
258 relationships, and this selectivity may make it better suited to detect genuine biologically  
259 relevant associations. Second, treeClust is exceptionally robust towards outlier measurements.  
260 Next, we tested which of these may be relevant for ProteomeHD.

## 261 **Outliers in ProteomeHD**

262 We first evaluated the impact of outlier measurements in ProteomeHD. We used two different  
263 statistical methods to detect outliers, as they identify distinct types of outliers. Data points with  
264 large studentized residuals are regression outliers, meaning they are far from the regression line  
265 but not necessarily unusual with regards to the overall distribution of the ratios (Fig. 3A). This  
266 type of outlier may lead to an underestimation of the real correlation coefficient. In contrast,  
267 outliers with a large Mahalanobis distance are far from the bulk of the data but can be close to  
268 the regression line (Fig. 3A). These outliers may lead to an overestimation of the real correlation  
269 coefficient. In principle, regression and Mahalanobis outliers could be seen as real-life examples  
270 of the outliers shown in Anscombe's third and fourth dataset, respectively.

271 We then tested if either of these outlier types explains why treeClust outperforms PCC  
272 for ProteomeHD data. For this we compared protein pairs that scored high (i.e. ranking in the  
273 top 0.1% pairs) with one method but were not detected (i.e. not ranked in top 0.5% pairs) by the  
274 other. This resulted in the following two groups of protein pairs: (a) 8,786 protein pairs with high  
275 treeClust similarities that were not detected as co-regulated by their PCCs; (b) 9,593 protein  
276 pairs with high PCCs that were not detected by treeClust. Functional annotation of these groups  
277 using a gold standard based on Reactome revealed that 60% of protein pairs found exclusively  
278 by treeClust are known true-positive associations, compared to only 13% of the PCC-specific  
279 pairs (Fig. 3B). Therefore, treeClust-specific co-regulation pairs are predominantly true  
280 interactions missed by PCC, whereas PCC-specific pairs are mostly unrelated proteins falsely  
281 identified as co-regulated by PCC. Similar distributions were observed when comparing  
282 treeClust to Spearman's rho and bicor (Fig. 3B).

283 Next, we asked if regression outliers in ProteomeHD may explain the difference between  
284 these groups. Surprisingly, we find that the number of regression outliers is very similar for  
285 treeClust-specific and PCC-specific protein pairs (5.8% vs 5.9% on average), as well as  
286 treeClust and rho-specific and bicor-specific pairs (Fig. 3C). However, the impact of outliers may  
287 not just stem from their number but also from their actual position and distribution compared to  
288 the rest of the data. We therefore removed outliers and measured the effect this had on the  
289 correlation coefficients. Indeed, removing regression outliers has a strong impact on  
290 co-regulation pairs that had been detected by treeClust but not by PCC, increasing their  
291 average PCC by 0.15 (Fig. 3D). This suggests that treeClust-specific co-regulation pairs tend to  
292 be genuine, biologically relevant interactions that are missed by PCC due to regression outliers.  
293 However, removing regression outliers had no dramatic effect on pairs that had been missed by  
294 rho or bicor (Fig. 3D).

295 In contrast to regression outliers, Mahalanobis outliers are clearly enriched among pairs  
296 only detected by PCC (27% vs 22% on average), or only by rho or bicor (Fig. 3E). Removing  
297 Mahalanobis outliers has a striking impact on PCC-specific pairs, reducing their average PCC  
298 by -0.29 (Fig. 3F). This indicates that co-regulated pairs detected by PCC - but not treeClust -  
299 are predominantly false-positive interactions whose high PCC is driven by Mahalanobis-type  
300 outliers. Associations detected only by rho or bicor, although enriched for Mahalanobis outliers,  
301 do not lose their high correlation coefficients by removing these outliers (Fig. 3F).



302 In summary, these results suggest that outliers are a key factor explaining why treeClust  
303 outperforms PCC in the analysis of ProteomeHD data. However, its improvement over rho and  
304 bicor is unlikely to be due to better outlier handling.

### 305 **Goodness-of-fit of genuine associations in ProteomeHD**

306 Increasing the scatter of values around the regression line led to a dramatic reduction of  
307 treeClust similarity in synthetic data (Fig. 2J). To quantify the overall difference between two  
308 proteins in real biological data we use the mean absolute error (MAE). Two protein pairs with  
309 very similar correlation coefficients can have vastly different MAEs (Fig. 4A). As expected, of the  
310 example pairs shown in Fig. 4A, the pair with the small MAE reflects a genuine biological  
311 association and receives a high treeClust similarity. In contrast, the pair with the large MAE is  
312 composed of two unrelated proteins and is not detected as co-regulated by treeClust (Fig. 4A).  
313 This suggests that treeClust may distinguish better than correlation measures between close,  
314 real interactions and loose, biologically irrelevant trends.

315 To assess this possibility in a systematic way we analysed the MAE distribution of all  
316 protein pairs that receive high correlation coefficients but low treeClust similarities, and vice  
317 versa. We find that protein pairs exclusively detected by rho or bicor tend to have much higher  
318 MAEs than those exclusively detected by treeClust (Fig. 4B). Interestingly, the difference in MAE  
319 distribution is not as pronounced between treeClust and PCC. Taken together this suggests that  
320 treeClust outperforms PCC mainly due to its outlier handling, whereas its improvement over rho  
321 and bicor is predominantly due to treeClust taking into account the “goodness-of-fit” of an  
322 association.

### 323 **Lack of non-linear relationships in ProteomeHD**

324 The selectivity of treeClust for linear relationships implies that it may fail to detect non-linear  
325 relationships that may be biologically relevant. We therefore investigated whether any genuine  
326 non-linear protein - protein associations exist in ProteomeHD. For this we fitted linear,  
327 exponential and logistic models to the correlation- or treeClust-specific protein pairs. For each  
328 pair we then select the best-fitting model based on the residual sum of squares (RSS). We find  
329 that exponential models rarely fit better than the linear regression models, but surprisingly,  
330 logistic models often do (Supplementary Fig. S1A). However, closer inspection of the data  
331 reveals that these cases are not genuine exponential or sigmoid relationships (Supplementary  
332 Fig. S1B). In contrast, the improved fit of the non-linear models is driven by Mahalanobis-type  
333 outliers. Removing these outliers also drastically decreases the number of instances in which  
334 non-linear models fit better than linear ones (Supplementary Fig. S1A). In summary, we have  
335 not been able to identify any clear non-linear relationships in ProteomeHD.

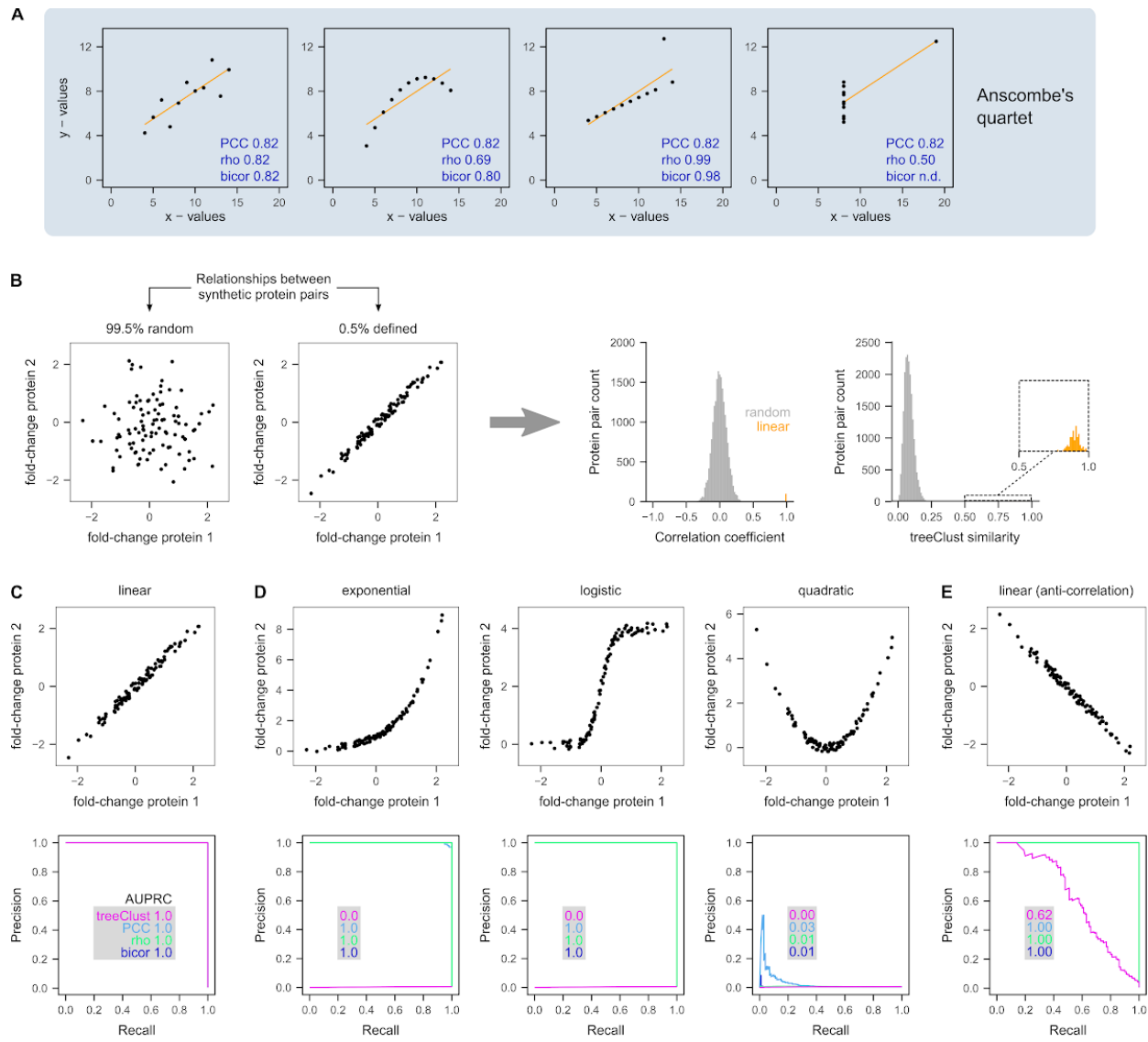
## 336 **CONCLUSION**

337 Having found treeClust to be a powerful alternative to correlation metrics for the detection of  
338 protein - protein links in proteomics data recently (13) we here demonstrated possible reasons

339 for this observation. treeClust is exceptionally robust against outliers and only identifies  
340 close-fitting, positive linear associations. In real proteomics datasets, these type of associations  
341 appear to be the biologically most relevant ones. Obvious disadvantages of using unsupervised  
342 machine-learning for this task are the required size and composition of the input data. At the  
343 moment, few proteomics datasets exists that are large enough, i.e. covering hundreds of  
344 conditions for thousands of genes, while maintaining a sufficiently low percentage of missing  
345 values for treeClust to work effectively. This applies to ProteomeHD and is likely going to  
346 become more prevalent in the near future, also thanks to the efforts of the ProteomeXchange  
347 consortium (35, 36). However, also smaller datasets can be analysed, by applying the algorithm  
348 many times to different subsets of the data and collecting the average similarities across these  
349 models (a bootstrapping approach) (13). It will be interesting to see how treeClust fares with  
350 other omics data types and other application areas where currently correlation approaches are  
351 in use.

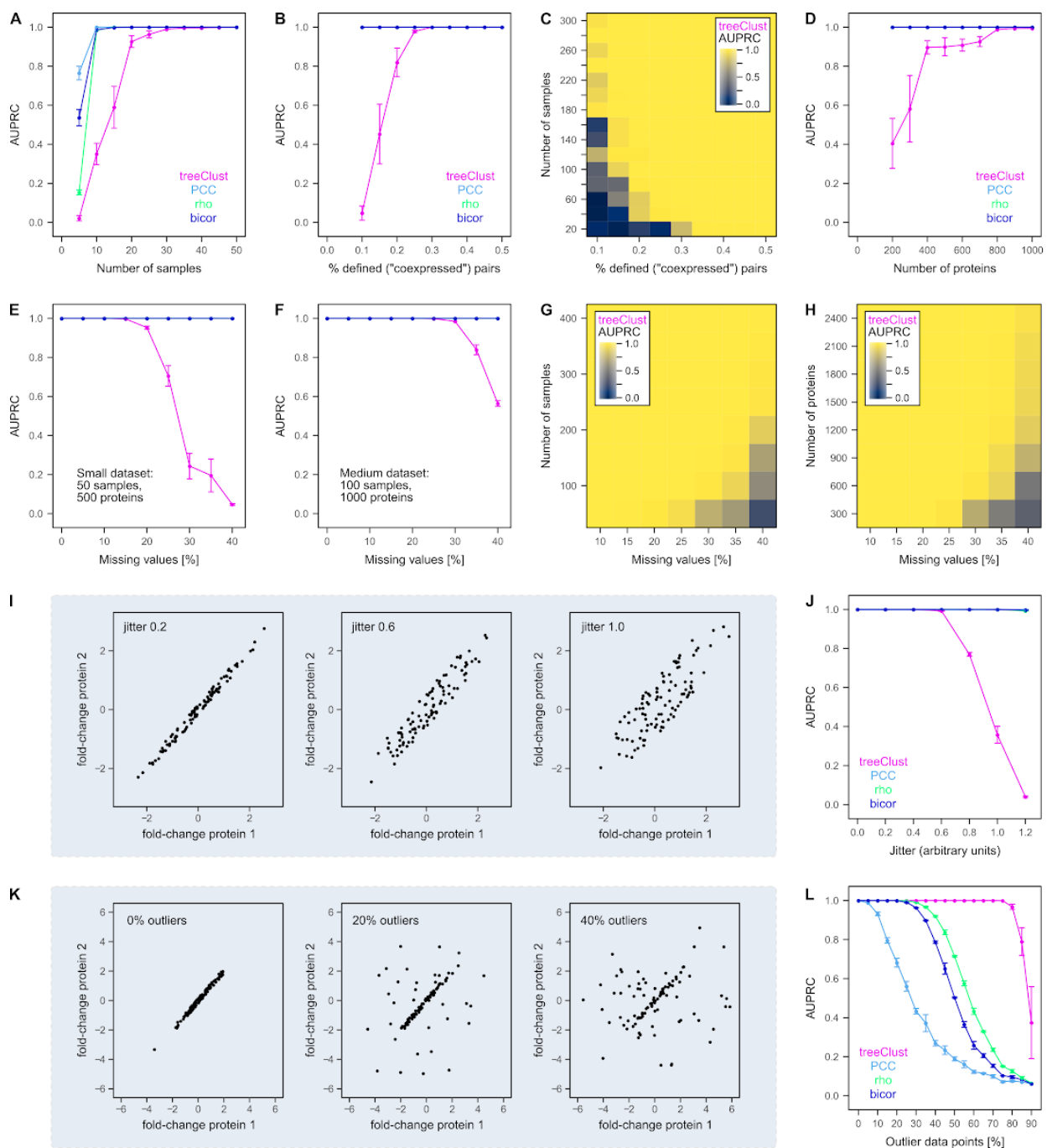
## 352 **ACKNOWLEDGEMENTS**

353 This work was supported by the Wellcome Trust through a Senior Research Fellowship to J.R.  
354 (grant number 103139). The Wellcome Centre for Cell Biology is supported by core funding from  
355 the Wellcome Trust (grant number 203149).



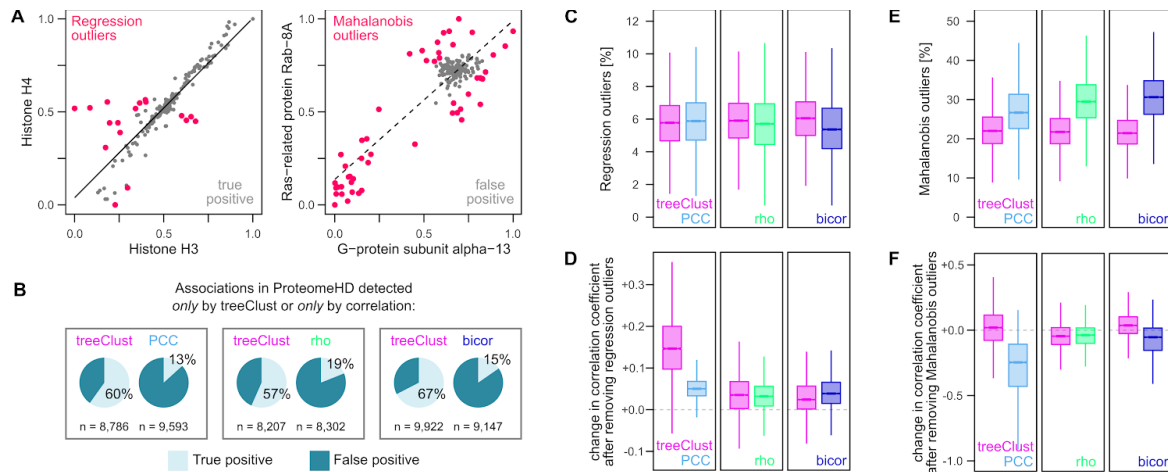
**Figure 1. treeClust detects specifically positive linear associations.**

(A) Anscombe's quartet is a collection of four 11-point datasets that have the same Pearson's correlation but very different relationships. If these variables were genes and the values were expression measurements, PCC would significantly underestimate (example 2 and 3) or overestimate (example 4) the extent of their coexpression. (B) To test how treeClust deals with such data we created a synthetic dataset consisting of 100 variables and 200 proteins. The dataset is designed such that out of all possible 19,900 combinations between these proteins, 0.5% have a defined relationship while the remaining 99.5% of pairs have not. Histograms show the resulting distribution of correlation coefficients or treeClust similarities. In this best-case scenario a complete separation between random and defined pairs is easily achieved. (C) Precision - recall (PR) analyses show that treeClust separates linear from random relationships perfectly, resulting in an area under the PR curve (AUPRC) of 1. The same result is observed for the three tested correlation-based metrics: PCC, Spearman's rho and biweight midcorrelation (bicor). The four PR curves overlap fully. (D) TreeClust completely fails to detect exponential or logistic relationships (AUPRC = 0). In contrast, these pairs still score high enough with PCC, rho and bicor to be completely separated from the pool of random associations. No metric detects quadratic relationships. (E) Anti-correlations are not identified well by treeClust.



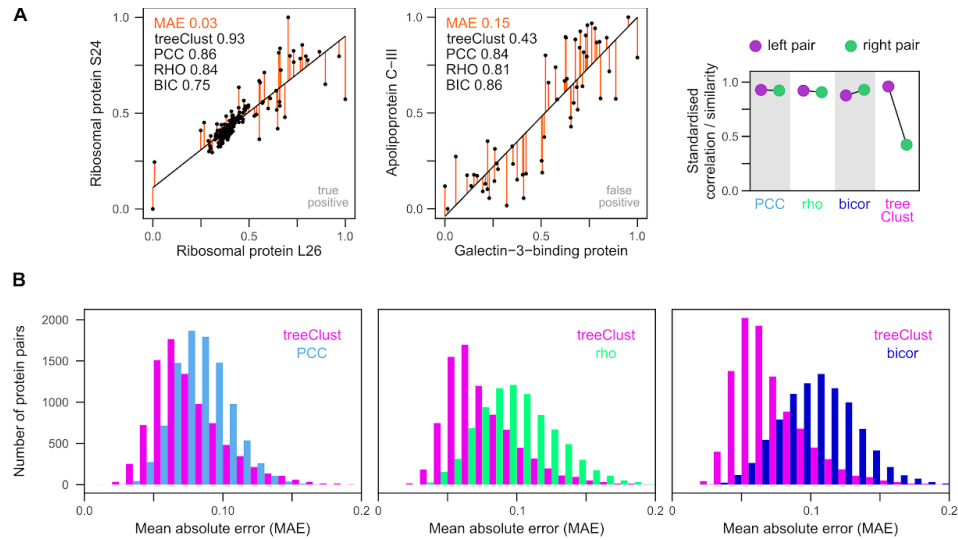
**Figure 2. Using synthetic data to benchmark treeClust performance**

(A) Increasing the number of samples (variables) in a synthetic dataset improves the performance of treeClust, Pearson's correlation (PCC), Spearman's correlation ( $\rho$ ) and Biweight Midcorrelation (bicor). TreeClust needs more samples for optimal performance than the correlation metrics. Each synthetic dataset, containing 500 proteins and 0.3% linear relationships, was created in triplicate. Points show the average area under the precision recall curve (AUPRC) obtained for each setting. Error bars show the standard error of the mean. (B) Same as A but increasing the linear associations in datasets of 50 samples and 500 proteins. This has no impact on the three correlation metrics, so their curves overlap fully at AUPRC 1. (C) Combinatorial impact of the two parameters on treeClust AUPRC is shown through a colour-gradient.  $N$  proteins = 500. (D) Same as A but increasing the number of observations (proteins).  $N$  samples = 20, 0.3% linear associations. (E, F) Adding missing values to a small and medium dataset, respectively. Points show the average area under the precision recall curve (AUPRC) on treeClust performance.  $N$  proteins = 1,000. (H) Same as G but modifying the number of proteins ( $N$  samples = 150). (I) Scatterplots illustrating the effect of increasing the difference between variables. (J) This increasing dispersion strongly impairs detection of linear associations by treeClust, but not correlation metrics. (K) Scatterplots illustrating the effect of adding outlier data points. (L) Impact of outlier data on the four coexpression measures. TreeClust is exceptionally robust against outliers.



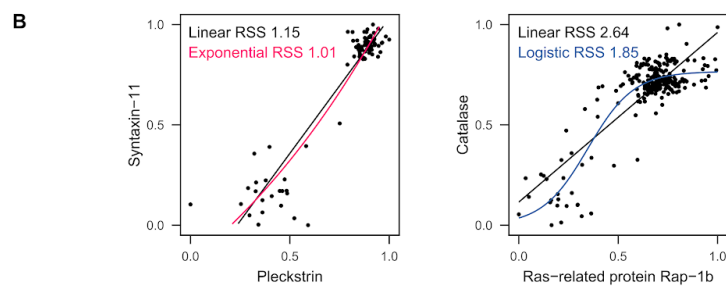
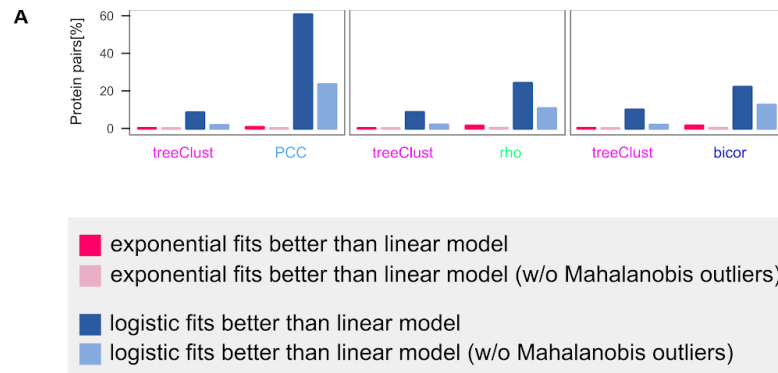
**Figure 3. Outliers in ProteomeHD and their impact on correlation metrics**

(A) Two examples protein pairs from ProteomeHD illustrate the two different types of outliers. Regression outliers are detected via studentized residuals and are located far away from the regression line. These outliers will decrease correlation coefficients. Outliers detected via their Mahalanobis distance are located far away from the bulk of the data, but can be close to the regression line. As in the example show, these outliers can cause high correlation coefficients even if no general correlation exists. Fold-changes have been scaled to lie between 0 and 1. (B) Co-regulated protein pairs were divided into those detected by treeClust but not by PCC and vice versa. Separate comparisons were made for pairs detected by treeClust but not by rho, and treeClust but not bicor. The pairs in the resulting groups were annotated using Reactome into known, biologically relevant interactions (true positives) and pairs that were unlikely to have any biological associations (false positives). Note that treeClust-specific pairs tend to be true positives, whereas correlation-specific pairs tend to be false positives. (C) The number of regression outliers is very similar in all six groups. (D) Removing regression outliers increases the PCC of protein pairs that were previously detected only by treeClust, suggesting PCC missed some of these pairs because of regression outliers. This is not the case for pairs missed by rho or bicor. (E) Mahalanobis outliers are more frequent in protein pairs detected by all three correlation metrics than pairs specific to treeClust. (F) Removing Mahalanobis outliers decreases the PCC of pairs that were originally detected only by PCC, suggesting their PCC was supported mainly by the Mahalanobis outliers. The correlation of pairs detected only by rho or bicor is not affected strongly by removing Mahalanobis outliers.



**Figure 4. Goodness-of-fit is a critical parameter for detecting genuine associations in ProteomeHD**

(A) Two examples pairs from ProteomeHD to illustrate close and loose - fitting regression. Both pairs have similar regression slopes and correlation coefficients. However, the left pair has a much smaller mean absolute error (MAE). MAE is the average of the absolute residuals (orange). The left pair is a known, biologically relevant interaction documented by the Reactome gold standard, the right pair is not. Unlike the correlation metrics, treeClust assigns the right pair are much weaker similarity. The difference between the left and right example pair is shown by the scatterplot on the right. For this plot the correlation metrics and treeClust similarities were standardised to fall within a range of [0,1] to make them comparable. (B) Systematic comparison of MAEs from protein pairs that scored high both with treeClust and PCC (or rho, or bicor), or pairs that scored high with either metric alone. Protein pairs exclusively detected by correlation metrics tend to have much higher MAEs, possibly explaining why they are predominantly false-positive hits.



### Supplementary Figure S1. Lack of genuine non-linear relationships in ProteomeHD

(A) Exponential and logistic (sigmoid) models were fitted to all protein pairs that scored high with treeClust or the three correlation metrics. Model fit was compared through their residual sum of squares (RSS). Exponential models only fitted better than linear ones in rare cases, but logistic models often did. Around half of the protein pairs detected specifically by PCC are better explained by a logistic than a linear model. However, this is mainly driven by Mahalanobis-type outliers. Removing those strongly reduces the number of logistic models outfitting the linear ones. (B) Two example regressions where an exponential (left) or logistic (right) model fits better than a linear one. Note that this clearly reflects overfitting due to outliers rather than genuine non-linear relationships.

356 **REFERENCES**

- 357
- 358 1. DeRisi, J. L., Iyer, V. R., and Brown, P. O. (1997) Exploring the metabolic and genetic control of gene expression on a genomic scale. *Science* 278, 680–686
- 359 2. Eisen, M. B., Spellman, P. T., Brown, P. O., and Botstein, D. (1998) Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. U. S. A.* 95, 14863–14868
- 360
- 361
- 362 3. Kim, S. K., Lund, J., Kiraly, M., Duke, K., Jiang, M., Stuart, J. M., Eizinger, A., Wylie, B. N., and Davidson, G. S. (2001) A gene expression map for *Caenorhabditis elegans*. *Science* 293, 2087–2092
- 363
- 364
- 365 4. Hughes, T. R., Marton, M. J., Jones, A. R., Roberts, C. J., Stoughton, R., Armour, C. D., Bennett, H. A., Coffey, E., Dai, H., He, Y. D., Kidd, M. J., King, A. M., Meyer, M. R., Slade, D., Lum, P. Y., Stepaniants, S. B., Shoemaker, D. D., Gachotte, D., Chakraburtt, K., Simon, J., Bard, M., and Friend, S. H. (2000) Functional discovery via a compendium of expression profiles. *Cell* 102, 109–126
- 366
- 367
- 368
- 369
- 370 5. Stuart, J. M., Segal, E., Koller, D., and Kim, S. K. (2003) A gene-coexpression network for global discovery of conserved genetic modules. *Science* 302, 249–255
- 371
- 372 6. van Dam, S., Vösa, U., van der Graaf, A., Franke, L., and de Magalhães, J. P. (2018) Gene co-expression analysis for functional classification and gene–disease predictions. *Brief. Bioinform.* 19, 575–592
- 373
- 374
- 375 7. Song, L., Langfelder, P., and Horvath, S. (2012) Comparison of co-expression measures: mutual information, correlation, and model based indices. *BMC Bioinformatics* 13, 328
- 376
- 377 8. Jaskowiak, P. A., Campello, R. J. G. B., and Costa, I. G. (2014) On the selection of appropriate distances for gene expression data clustering. *BMC Bioinformatics* 15 Suppl 2, S2
- 378
- 379
- 380 9. Wang, J., Ma, Z., Carr, S. A., Mertins, P., Zhang, H., Zhang, Z., Chan, D. W., Ellis, M. J. C., Townsend, R. R., Smith, R. D., McDermott, J. E., Chen, X., Paulovich, A. G., Boja, E. S., Mesri, M., Kinsinger, C. R., Rodriguez, H., Rodland, K. D., Liebler, D. C., and Zhang, B. (2017) Proteome Profiling Outperforms Transcriptome Profiling for Coexpression Based Gene Function Prediction. *Mol. Cell. Proteomics* 16, 121–134
- 381
- 382
- 383
- 384
- 385 10. Lapek, J. D., Jr, Greninger, P., Morris, R., Amzallag, A., Pruteanu-Malinici, I., Benes, C. H., and Haas, W. (2017) Detection of dysregulated protein–association networks by high-throughput proteomics predicts cancer vulnerabilities. *Nat. Biotechnol.* 35, 983–989
- 386
- 387



- 388 11. Kustatscher, G., Grabowski, P., and Rappsilber, J. (2017) Pervasive coexpression of  
389 spatially proximal genes is buffered at the protein level. *Mol. Syst. Biol.* 13, 937
- 390 12. Grabowski, P., Kustatscher, G., and Rappsilber, J. (2018) Epigenetic Variability Confounds  
391 Transcriptome but Not Proteome Profiling for Coexpression-based Gene Function  
392 Prediction. *Mol. Cell. Proteomics* 17, 2082–2090
- 393 13. Kustatscher, G., Grabowski, P., Schrader, T., Passmore, J. B., Schrader, M., and  
394 Rappsilber, J. (2019) The human proteome co-regulation map reveals functional  
395 relationships between proteins. *bioRxiv* 582247; doi: <https://doi.org/10.1101/582247>
- 396 14. Buttrey, S. E., and Whitaker, L. R. (2015) treeClust: an R package for tree-based clustering  
397 dissimilarities. *The R Journal* 7, 227–236
- 398 15. Buttrey, S. E., and Whitaker, L. R. (2016) A scale-independent, noise-resistant dissimilarity  
399 for tree-based clustering of mixed data. *NPS Technical Report Archive*,
- 400 16. Therneau, T. M. (1983) A short introduction to recursive partitioning. *Orion Technical Report*  
401 21,
- 402 17. Therneau, T. M., Atkinson, E. J., and Others (1997) An introduction to recursive partitioning  
403 using the RPART routines. *Technical Report 61, Mayo Clinic*,
- 404 18. Breiman, L. (2001) Random Forests. *Mach. Learn.* 45, 5–32
- 405 19. Stekhoven, D. J., and Bühlmann, P. (2012) MissForest--non-parametric missing value  
406 imputation for mixed-type data. *Bioinformatics* 28, 112–118
- 407 20. Gower, J. C. (1971) A General Coefficient of Similarity and Some of Its Properties.  
408 *Biometrics* 27, 857–871
- 409 21. Zoratti, E. M., Krouse, R. Z., Babineau, D. C., Pongracic, J. A., O'Connor, G. T., Wood, R.  
410 A., Khurana Hershey, G. K., Kerckmar, C. M., Gruchalla, R. S., Kattan, M., Teach, S. J.,  
411 Sigelman, S. M., Gergen, P. J., Togias, A., Visness, C. M., Busse, W. W., and Liu, A. H.  
412 (2016) Asthma phenotypes in inner-city children. *J. Allergy Clin. Immunol.* 138, 1016–1029
- 413 22. Oliveira, S., Zêzere, J. L., Queirós, M., and Pereira, J. M. (2017) Assessing the social  
414 context of wildfire-affected areas. The case of mainland Portugal. *Appl. Geogr.* 88, 104–117
- 415 23. R Core Team (2018) R: A Language and Environment for Statistical Computing.
- 416 24. Dowle, M., and Srinivasan, A. (2018) data.table: Extension of `data.frame`.
- 417 25. Wickham, H. (2016) *ggplot2: Elegant Graphics for Data Analysis* (Springer)

- 418 26. Auguie, B. (2017) gridExtra: Miscellaneous Functions for “Grid” Graphics.
- 419 27. Wilke, C. O. (2018) cowplot: Streamlined Plot Theme and Plot Annotations for “ggplot2.”
- 420 28. Garnier, S. (2018) viridis: Default Color Maps from “matplotlib.”
- 421 29. Ong, S.-E., Blagoev, B., Kratchmarova, I., Kristensen, D. B., Steen, H., Pandey, A., and  
422 Mann, M. (2002) Stable isotope labeling by amino acids in cell culture, SILAC, as a simple  
423 and accurate approach to expression proteomics. *Mol. Cell. Proteomics* 1, 376–386
- 424 30. Fabregat, A., Sidiropoulos, K., Garapati, P., Gillespie, M., Hausmann, K., Haw, R., Jassal,  
425 B., Jupe, S., Korninger, F., McKay, S., Matthews, L., May, B., Milacic, M., Rothfels, K.,  
426 Shamovsky, V., Webber, M., Weiser, J., Williams, M., Wu, G., Stein, L., Hermjakob, H., and  
427 D’Eustachio, P. (2016) The Reactome pathway Knowledgebase. *Nucleic Acids Res.* 44,  
428 D481–7
- 429 31. Langfelder, P., and Horvath, S. (2008) WGCNA: an R package for weighted correlation  
430 network analysis. *BMC Bioinformatics* 9, 559
- 431 32. Langfelder, P., and Horvath, S. (2012) Fast R Functions for Robust Correlations and  
432 Hierarchical Clustering. *J. Stat. Softw.* 46,
- 433 33. Grau, J., Grosse, I., and Keilwagen, J. (2015) PRROC: computing and visualizing  
434 precision-recall and receiver operating characteristic curves in R. *Bioinformatics* 31,  
435 2595–2597
- 436 34. Anscombe, F. J. (1973) Graphs in Statistical Analysis. *Am. Stat.* 27, 17–21
- 437 35. Vizcaíno, J. A., Deutsch, E. W., Wang, R., Csordas, A., Reisinger, F., Ríos, D., Dienes, J.  
438 A., Sun, Z., Farrah, T., Bandeira, N., Binz, P.-A., Xenarios, I., Eisenacher, M., Mayer, G.,  
439 Gatto, L., Campos, A., Chalkley, R. J., Kraus, H.-J., Albar, J. P., Martinez-Bartolomé, S.,  
440 Apweiler, R., Omenn, G. S., Martens, L., Jones, A. R., and Hermjakob, H. (2014)  
441 ProteomeXchange provides globally coordinated proteomics data submission and  
442 dissemination. *Nat. Biotechnol.* 32, 223–226
- 443 36. Vizcaíno, J. A., Csordas, A., del-Toro, N., Dienes, J. A., Griss, J., Lavidas, I., Mayer, G.,  
444 Perez-Riverol, Y., Reisinger, F., Ternent, T., Xu, Q.-W., Wang, R., and Hermjakob, H. (2016)  
445 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res.* 44, D447–56