

Denoising large-scale biological data using network filters

Andrew J Kavran

University of Colorado Boulder <https://orcid.org/0000-0001-7660-3922>

Aaron Clauset (✉ aaron.clauset@colorado.edu)

University of Colorado Boulder <https://orcid.org/0000-0002-3529-8746>

Research article

Keywords: Networks, Denoising, Machine Learning

Posted Date: January 14th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-66071/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on March 25th, 2021. See the published version at <https://doi.org/10.1186/s12859-021-04075-x>.

RESEARCH

Denoising large-scale biological data using network filters

Andrew J Kavran^{1,2} and Aaron Clauset^{3,2,4*}

*Correspondence:

aaron.clauset@colorado.edu

¹Department of Biochemistry,
University of Colorado, Boulder,
CO, USA

Full list of author information is
available at the end of the article

Abstract

Background: Large-scale biological data sets are often contaminated by noise, which can impede accurate inferences about underlying processes. Such measurement noise can arise from endogenous biological factors like cell cycle and life history variation, and from exogenous technical factors like sample preparation and instrument variation.

Results: We describe a general method for automatically reducing noise in large-scale biological data sets. This method uses an interaction network to identify groups of correlated or anti-correlated measurements that can be combined or “filtered” to better recover an underlying biological signal. Similar to the process of denoising an image, a single network filter may be applied to an entire system, or the system may be first decomposed into distinct modules and a different filter applied to each. Applied to synthetic data with known network structure and signal, network filters accurately reduce noise across a wide range of noise levels and structures. Applied to a machine learning task of predicting changes in human protein expression in healthy and cancerous tissues, network filtering prior to training increases accuracy up to 43% compared to using unfiltered data.

Conclusions: Network filters are a general way to denoise biological data and can account for both correlation and anti-correlation between different measurements. Furthermore, we find that partitioning a network prior to filtering can significantly reduce errors in networks with heterogeneous data and correlation patterns, and this approach outperforms existing diffusion based methods. Our results on proteomics data indicate the broad potential utility of network filters to applications in systems biology.

Keywords: Networks; Denoising; Machine Learning

¹Background

²System-wide **molecular profiling data** are often contaminated by noise, which can²
³obscure biological signals of interest. Such noise can arise from both endogenous³
⁴biological factors and exogenous technical factors. **These** factors include reagent⁴
⁵and protocol variability, researcher technique, passage number effects, stochastic⁵
⁶gene expression, and cell cycle asynchronicity. This variability can mask underlying⁶
⁷biological signals when measuring cell state and how it changes under different⁷
⁸conditions, e.g., in development [1, 2], cancer progression [3], and adaptive drug⁸
⁹resistance [4, 5]. Noise has also been implicated in the appearance of false signals⁹
¹⁰and in the non-replicability of some studies [6, 7]. Identifying and correcting noisy¹⁰
¹¹measurements before analysis is likely to improve the detection of subtle biological¹¹
¹²signals and enable more accurate predictions in systems biology. 12

¹³If correlations between related molecular signals are stronger than correlations 13
¹⁴among sources of noise, then distinct but related signals can be combined to de- 14
¹⁵noise biological measurements, at the expense of a smaller effective sample size. 15
¹⁶There are three common approaches to identifying related signals: gene sets, sub- 16
¹⁷space embedding, and networks. In the first category, methods like GSEA [8, 9] 17
¹⁸use the enrichment of genes within curated sets to project the data onto biologi- 18
¹⁹cally relevant features. While gene sets can increase the power to identify differen- 19
²⁰tially regulated processes, they are inherently coarse, and can themselves be noisy, 20
²¹incomplete, or biased, and thus may not generalize to novel processes. Subspace 21
²²embedding techniques include PCA [10], clustering [11], and neural network au- 22
²³toencoders [12, 13]. These methods can capture novel gene-gene correlations, but 23
²⁴they rarely incorporate biological information into the feature extraction, which can 24
²⁵limit both interpretability and generalizability. 25

²⁶Molecular profiling data alone does not directly inform which measurements 26
²⁷should be more or less related to each other. Networks that represent a molecular 27
²⁸system's functional structure can provide this missing information. For example, 28
²⁹protein-protein interaction, metabolic reaction, and gene regulation networks each 29
³⁰encode precise and biologically meaningful information about which groups of mea- 30
³¹sured protein expression levels, metabolite concentrations, or transcript levels are 31
³²functionally related, and hence which measurements should be combined to filter 32
³³out independent noise. Current network approaches use computationally intensive 33

1 methods to identify which entities are most related, which can limit their utility for
2 large networks and general usability [14, 15]

3 Among neighboring elements in the network, the underlying signals may be cor-
4 related (assortative) or anti-correlated (disassortative) [16]. For example, differen-
5 tial expression tends to correlate between neighboring genes in a regulatory net-
6 work [17]. In contrast, inhibitory or compensatory interactions [18, 19] will tend
7 to produce a disassortative relationship. Beyond pairs of measurements, networks
8 can also exhibit large-scale mixing patterns among these interactions, such that
9 a network may be more or less assortative in some regions and disassortative in
10 others [20]. Existing network-based methods typically do not exploit this variabil-
11 ity, and instead assume globally assortative mixing by applying a single filter to the
12 whole network [21, 14, 15]. Mismatching the filter and the relationship type, e.g., an
13 assortative filter with anti-correlated measurements, can further obscure the under-
14 lying biological signals. Here, we describe a general network-based method that can
15 automatically detect large-scale mixing patterns and account for both assortative
16 and disassortative relationships.

17 These network filters are closely related to kernel-based methods in image pro-
18 cessing [22], in which groups of related pixels are transformed together to improve
19 their underlying visual signal. Most such techniques leverage an image’s underlying
20 grid geometry to choose which pixels have related signals for denoising. Networks
21 lack this geometry because a node’s interactions are inherently unordered, whereas
22 the left- and right-hand neighbors of a pixel are clearly defined. This connection
23 between network filters and image processing is rich with potentially useful ideas
24 that could be adapted to process large-scale biological data. For instance, commu-
25 nity detection in networks is a clear analog of the common “segmentation” step in
26 image analysis, in which pixels are first partitioned into groups that represent the
27 large-scale structure of an image, e.g., to separate foreground and background, or a
28 car from the street, and then different filters are applied to each segment (module).

29 We first describe two classes of network filters, which combine measurement val-
30 ues from neighboring nodes to calculate an assortative or disassortative denoised
31 value, and we describe a general algorithm that decomposes the network into struc-
32 tural modules and then automatically applies the most appropriate filter to the
33 nodes and connections within each module. When applied to synthetic data where

¹the true values and network structure are known, these filters substantially reduce¹
²errors relative to a baseline. In addition, we show how applying the wrong filter²
³with respect to the underlying biological relationship can lead to increased errors.³
⁴Finally, to test the practical utility of these methods in a more realistic setting,⁴
⁵we investigate the impact of network filtering on a machine learning task in which⁵
⁶we predict changes in human protein expression data when a healthy tissue be⁶
⁷comes cancerous. Using the network filters to denoise the expression data before⁷
⁸model training increases the subsequent prediction accuracy up to 43% compared⁸
⁹to training on unfiltered data. 9

10

10

11

11

¹²Results 12

¹³Network filters 13

14

14

¹⁵A network filter is specified by a function $f[i, \mathbf{x}, G]$, which takes as input the index¹⁵
¹⁶of the measurement (node) to be denoised, the list of all measurements \mathbf{x} , and the¹⁶
¹⁷network structure G among those measurements. The output is the denoised value¹⁷
¹⁸ \hat{x}_i . Here, we consider only local network filters, which use the measurement values¹⁸
¹⁹of i 's immediate neighbors in G , denoted by the node set ν_i , which are likely to be¹⁹
²⁰the most biologically relevant for denoising. Each filter is applied synchronously, so²⁰
²¹that all denoised values are obtained simultaneously to prevent feedback within the²¹
²²denoising process. 22

²³We note that the idea of a network filter can naturally generalize to exploit infor-²³
²⁴mation, if available, about the sign or strength of interactions in G . This information²⁴
²⁵can be encoded by an edge weight w_{ij} , which can capture inhibitory or excitatory²⁵
²⁶interactions that are strong or weak. Below, we focus on the case in which this²⁶
²⁷information is not available. 27

²⁸When a measurement x_i correlates with the values of its neighbors x_{ν_i} in the ²⁸
²⁹network (assortativity), a network filter should adjust x_i to be more similar to the ²⁹
³⁰measured values of its neighbors (Fig. 1A). Among the many choices of functions ³⁰
³¹with this qualitative behavior, the mean and median have useful mathematical ³¹
³²properties, and connect with past work [21]. This setting is analogous to a smoothing ³²
³³operation in image processing, in which a pixel's value is replaced by the mean or ³³

1 median of its value and its neighbors' values. In the context of a network, the mean
 2 and median “smoothing” filters have the forms:

$$3 \quad f_{\bullet,1}[i, \mathbf{x}, G] = \frac{1}{1 + k_i} \left(x_i + \sum_{j \in \nu_i} x_j w_{ij} \right), \quad (1) 5$$

6 where $w_{ij} = 1$ and k_i is the degree of node i , reflecting unweighted interactions, and
 7

$$8 \quad f_{\bullet,2}[i, \mathbf{x}, G] = \text{median}\{x_i, x_{\nu_i} w_{ij}\}. \quad (2) 9$$

10
 11 When a measurement x_i anti-correlates with the values of its neighboring nodes,
 12 a network filter should adjust x_i to be more distant from its neighbors (Fig. 1A).
 13 This setting is analogous to enhancing the contrast in an image, e.g., when using
 14 the technique of unsharp masking to enhance the high frequency signal in an image
 15 to make it sharper. In the context of a network, this “sharpening” filter has the
 16 form:

$$17 \quad f_{\circ}[i, \mathbf{x}, G] = \alpha(x_i - f_{\bullet,1}[i, \mathbf{x}, G]) + \bar{x} \quad (3) 18$$

19 where α is a constant scaling factor, and $\bar{x} = n^{-1} \sum_i x_i$ is the global mean. Because
 20 α is a free parameter, its value should be determined de novo for each data set. For
 21 the data sets in this study, we empirically determined the optimal $\alpha = 0.8$ using
 22 cross validation.

23 When a system exhibits large-scale mixing patterns of assortative and disassor-
 24 tative relationships, a network should first be partitioned into structural modules
 25 using a community detection algorithm, so that relationships within each module
 26 are more homogeneous. Let $\vec{s} = \mathcal{A}(G)$ denote the result of applying a community
 27 detection algorithm \mathcal{A} to network G , and say that G_{s_i} denotes the subgraph of
 28 nodes and connections within the module s_i that contains node i . Given such a
 29 modular decomposition \vec{s} , a filter can then be applied to only the subgraph G_{s_i}
 30 that contains measurement i . As a result, relationships that span the boundary
 31 between two modules will have no influence on the filtered values (Fig. 1B).

32 After partitioning, the same filter can be applied to every community, or sharp
 33 and smooth filters can be applied to communities with more or less assortative

values, respectively. We define such a “patchwork filter” as:

$$f[i, \mathbf{x}, G_{s_i}] = \begin{cases} f_{\circ}[i, \mathbf{x}, G_{s_i}] & r_{s_i} < 0 \\ f_{\bullet,1}[i, \mathbf{x}, G_{s_i}] & r_{s_i} \geq 0 \end{cases}, \quad (4)$$

where r_{s_i} is the standard assortativity coefficient calculated over observed values within community s_i [16]. While any community detection algorithm can be used for \mathcal{A} , here we use methods from three classes of algorithms: modularity maximization [23], spectral partitioning [24], and statistical inference. For community detection by statistical inference, we use the degree-corrected stochastic block model or DC-SBM [25] or the “metadata-aware” version of DC-SBM [26], which are considered state-of-the-art methods [27].

13 Tests using synthetic data

14 We evaluated the performance of these network filters in two controlled experiments
15 with either non-modular or modular synthetic networks, and varying structures and
16 levels of noise. And, we compare the performance of network filters to other network-
17 based denoising methods that combine values of nodes weighted by a diffusion
18 matrix [14, 15].

19 In the first experiment, we generated simple random graphs with heavy-tailed
20 degree distributions (see Methods) and assigned each node a value drawn from a
21 Normal distribution with mean $\mu = 100$ and standard deviation $\sigma = 10$. These val-
22 ues were drawn in such a way that the assortativity coefficient of the network ranged
23 from $r \in [-0.8, 0.8]$ (see Methods). As a result, connected values ranged from being
24 highly anticorrelated to highly correlated. To simulate independent measurement
25 noise, we permuted the values among a uniformly random 25% of nodes, and then
26 denoised these “corrupted” values. We find qualitatively similar results for other
27 choices of the fraction permuted. Results report the mean absolute error (MAE) of
28 a denoised value, averaged over 5000 replications.

29 Without a filter, the average error of a “denoised” value is independent of the
30 underlying correlation (assortativity) among connected values, because this nearby
31 information is left unexploited (Fig. 2A). In contrast, applying a network filter to
32 denoise the corrupted values can substantially improve their accuracy, depending on
33 how strongly coupled a measurement’s true value is with its neighbors’, and what

¹filter is applied to recover that information. For the particular parameters of this¹
²experiment, filtering can reduce the error by 37–50% over no filter, and by roughly²
³20% even in the case of uncorrelated signals ($r = 0$), due to a regression to the mean³
⁴effect. Error reductions are largest when a network “smoothing” filter is applied to⁴
⁵strongly assortative signals, and when a network “sharpening” filter is applied to⁵
⁶strongly disassortative signals. That is, denoising works best when the underlying⁶
⁷signal structure is matched with the assumptions of the filter. 7

⁸ When the wrong filter is applied, however, error rates can increase relative to not⁸
⁹filtering. In such a case, the filter creates more errors in the data than it corrects.⁹
¹⁰On the other hand, this “mismatch” penalty only degrades the overall accuracy at¹⁰
¹¹very high levels of correlation (anti-correlation) among the signals, where its magni-¹¹
¹²tude exceeds the natural benefits of filtering at all (Fig. 2A). When the underlying¹²
¹³correlations are moderate ($|r| < 0.4$), the average benefits of network filtering will¹³
¹⁴tend to outweigh the average error induced applying the wrong filter. 14

¹⁵ We also applied two other network methods to these non-modular synthetic¹⁵
¹⁶graphs. These methods denoise data by combining node values weighted by a dif-¹⁶
¹⁷fusion kernel. In the method called netSmooth [14], every node is weighted by a¹⁷
¹⁸personalized PageRank random walk vector [28], which are linearly combined to¹⁸
¹⁹create a denoised value. The second method is conceptually similar, but uses a¹⁹
²⁰graph Laplacian exponential diffusion kernel to weight nodes before linearly com-²⁰
²¹bining them to create the new denoised value [15]. Both methods have an adjustable²¹
²²parameter that determines the smoothness of the resulting denoised values. For²²
²³greater values of this smoothing parameter, the methods place less weight on the²³
²⁴node’s original noisy value and more weight on distant nodes. 24

²⁵ We applied both methods to the same synthetic random graphs as the network²⁵
²⁶filters, while varying the smoothing parameters between low smoothing (param-²⁶
²⁷eter = 0.1), and high smoothing (parameter = 0.9). The Laplacian exponential²⁷
²⁸kernel (Fig. 2B) and netSmooth (Fig. 2C) decrease the error of the noisy data²⁸
²⁹as the assortativity increases. Furthermore, both methods show lower error as the²⁹
³⁰smoothing parameter increases. These diffusion-based methods perform better than³⁰
³¹the smoothing network filters at either highly disassortative or weakly assortative³¹
³²values. Since these methods will typically use a larger number of node’s values to³²
³³denoise, their regression to the mean effect tends to be more accurate than the³³

¹more localized smoothing network filters. However, when a node's value becomes¹
²more correlated with its neighbors' values, the smoothing network filters decrease²
³the noise more than the diffusion-based methods. And while the diffusion based³
⁴methods work better than the smoothing filters on disassortative data, the sharp⁴
⁵filter is the best performing method for weakly to strongly disassortative data. ⁵

⁶ These tests assume that the network structure itself is not noisy. However, in real⁶
⁷ biological networks, there can be both missing and spurious edges [29]. We tested⁷
⁸ the robustness of network filters to noise in network structure. After creating a⁸
⁹ synthetic graph and assigning data to nodes, we add different levels of noise by re-⁹
¹⁰ placing true edges with new edges between nodes chosen uniformly at random [30].¹⁰
¹¹ Thus, this process simulates both cases where the network is missing edges and¹¹
¹² contains false edges. We find that a noisy network decreases the performance of the¹²
¹³ mean filter (Fig. S1A) and median filter (Fig. S1B) on graphs with assortative data¹³
¹⁴ ($r > 0$), and the sharp filter on graphs with disassortative data ($r < 0$) (Fig. S1C).¹⁴
¹⁵ However, the network filters still substantially reduce the error compared to the¹⁵
¹⁶ no-filter baseline. When the network is very noisy (90% rewired edges), applying¹⁶
¹⁷ a filter reduces the error compared to the no-filter baseline. This pattern is due to¹⁷
¹⁸ a regression to the mean effect, since rewiring the network effectively shrinks the¹⁸
¹⁹ assortativity coefficient closer to zero. ¹⁹

²⁰ In the second experiment, we again generated simple random graphs with heavy-²⁰
²¹ tailed degree distributions, but now also with modular structure, which better cap-²¹
²² tures the structure of empirical biological networks (see Methods). These modules²²
²³ denote groups of nodes that connect to other groups in statistically similar ways.²³
²⁴ For instance, protein interaction networks can be decomposed into groups with sim-²⁴
²⁵ ilar biological function, and these groups can have distinct types or levels of signal²⁵
²⁶ assortativity [20]. In this situation, applying a single filter to all parts of the net-²⁶
²⁷ work could introduce bias in the denoised values, by pooling nearby measurements²⁷
²⁸ indiscriminately, compared to filtering modules independently. ²⁸

²⁹ Here, we plant $\kappa = 5$ modules in the same kind of synthetic network as our first²⁹
³⁰ experiment, set each module to have a different mean value, and then vary the³⁰
³¹ fraction of modules that have a positive assortativity coefficient $|r| \in [0.4, 0.7]$ vs. a³¹
³² negative coefficient (see Methods). This kind of signal heterogeneity across modules³²
³³ mitigates the denoising benefits of a simple regression to the mean, and provides a³³

¹harder test for denoising methods. Given these choices, we generated values within¹
²a module, and simulated measurement noise as in the previous experiment (see²
³Methods). In addition to the previous filters, we also apply the “patchwork” filter³
⁴in this experiment. 4

⁵ As before, the average error of a denoised value with no filter provides a consistent⁵
⁶baseline against which we may assess improvements from filtering (Fig. 2D). And⁶
⁷similarly, the error for both the smooth and median filters falls steadily as the⁷
⁸fraction of modules with assortative signals increases. For the particular parameters⁸
⁹of this experiment, the median filter performs roughly 20% better than the mean⁹
¹⁰filter, reflecting the median’s well-known robustness to outliers, which arise here¹⁰
¹¹from the planted signal heterogeneity. 11

¹² The global sharp filter works poorly for all ratios when applied uniformly across¹²
¹³the whole network (Fig. S2). Because each module has a distinct mean value, the¹³
¹⁴global sharp filter generates errors by assuming the global mean is a good represen-¹⁴
¹⁵tation of the whole network. 15

¹⁶ In contrast, the patchwork filter with different community detection algorithms¹⁶
¹⁷exhibits less dynamic range in its error (Fig. 2D). When paired with the DC-SBM, it¹⁷
¹⁸is substantially more accurate than any other filter across different degrees of mod-¹⁸
¹⁹ular assortativity. For the particular parameters of this experiment, the patchwork¹⁹
²⁰filter paired with the DC-SBM reduces the mean error by 30–41% compared to no²⁰
²¹filtering and by 3–36% compared to median or mean filtering. Only when all of the²¹
²²modules are assortative does the median filter come close to the DC-SBM patch-²²
²³work filter’s accuracy. This advantage arises because the patchwork filter avoids²³
²⁴applying the same filter to different types of underlying signals, if the structure of²⁴
²⁵those signals correlates with the structure of the network (as it does here). That is,²⁵
²⁶applying a single filter to a modular network can introduce errors when denoising,²⁶
²⁷if the local mixing patterns across modules are heterogeneous. Pairing a commu-²⁷
²⁸nity detection algorithm with network filters can avoid this problem by identifying²⁸
²⁹large groups of nodes that should be filtered together, in much the same way that²⁹
³⁰different image filters can be applied after first segmenting an image into distinct³⁰
³¹regions. 31

³² However, for the modularity maximization and spectral partitioning algorithms,³²
³³the patchwork filter does not perform as well as when paired with the DC-SBM³³

1 because the algorithms do not partition the network as closely to the true com-
2 munity structure. Thus, the patchwork filter uses measurements from outside a
3 single community more often with these algorithms. Despite imperfect partitioning,
4 the patchwork filter paired with modularity and spectral partitioning algorithms
5 performs 14%-28% better than the mean filter across all levels of modular assorta-
6 tivity. The median filter outperforms the spectral patchwork (9%) and modularity
7 patchwork (15%) at the highest level of modular assortativity, but the patchwork
8 filter still outperforms, or matches the median filter across the rest of the levels of
9 modular assortativity.

10 We also applied the diffusion-based methods to these synthetic modular net-
11 works. The error for both the Laplacian exponential kernel (Fig. 2E) and netSmooth
12 (Fig. 2F) only slightly decreases as the fraction of modules with assortative signals
13 increases. In contrast to the non-modular case, increasing the smoothing paramete-
14 ter for both methods increases the error across all settings. This loss of accuracy
15 occurs because increasing the smoothing parameter places greater weight on more
16 distant nodes which are more likely to be drawn from a different distribution. Hence,
17 the diffusion kernels are more likely to combine values from nodes from different
18 communities leading to a higher error rate.

20 Denoising protein expression levels in cancer

21 To evaluate the utility of network filters for denoising biological data in realis-
22 tic settings, we construct a machine learning task in which we predict the precise
23 changes in human protein expression levels when a healthy tissue becomes cancer-
24 ous (see Methods). This task has potential applications to detecting pre-cancerous
25 lesions [31, 32]. We then quantify the improvement in out-of-sample prediction accu-
26 racy when using a network filter to denoise the input expression data before model
27 training, compared to training on unfiltered data.

28 For this experiment, protein expression data are drawn from the Human Protein
29 Atlas (HPA) [33], which provides large-scale immunohistochemistry (IHC) mea-
30 surements for over 12,000 human proteins in 20 tissues, each in both healthy and
31 cancerous states. Antibody based methods like IHC are known to be noisy and
32 prone to variation from uncontrolled experimental parameters [34], which makes
33 this data set a realistic example of noisy molecular profiling data. A standard prin-

1 principal component analysis (PCA) of the raw HPA expression data reveals that the¹
2 first component correlates with variations in tissue type, while the second correlates²
3 with differences between tissue state (healthy vs. cancerous) (Fig. 3A). Some tis-³
4 sues, however, change more than others, and the changes are not always in the same⁴
5 direction. Hence, predicting the precise changes represents a useful and non-trivial⁵
6 machine learning task that network filtering may improve. ⁶

7 For the network filters and diffusion-based methods, we use a comprehensive map⁷
8 of the human protein-protein interaction network (PPIN) [35], which combines data⁸
9 from several interactome databases and is curated for biological interactions with⁹
10 high levels of evidence. While this network represents a broad collection of authori-¹⁰
11 tative interactome data, the completeness of the human PPIN is still uncertain [29],¹¹
12 and we do not regard this network as itself noise-free. Taking the intersection of¹²
13 proteins contained in both expression data and interaction network (see Methods)¹³
14 yields data on $n = 8,199$ proteins in a network with $m = 37,607$ edges. ¹⁴

15 In the machine learning task, we perform a K -nearest neighbor regression on an¹⁵
16 embedded representation of the protein expression data to learn how expression¹⁶
17 levels change with tissue state (see Methods). We evaluate the trained model via¹⁷
18 the MAE between the predicted and the actual changes in protein expression under¹⁸
19 leave-one-out cross validation (in which we train on 19 tissue pairs, and predict¹⁹
20 on the 20th) with or without denoising the expression data with a network filter²⁰
21 or diffusion-based method prior to model training. Because the number K is a²¹
22 free parameter that controls the complexity of the learned model, we evaluate the²²
23 robustness of our results by systematically varying K . For the patchwork filter,²³
24 we partitioned the graph into 10 modules using the DC-SBM [25] or spectral algo-²⁴
25 rithm [24], while the modularity maximization algorithm [23] automatically chooses²⁵
26 the number of modules. Then, we apply the mean filter within each module. In this²⁶
27 data, most measured values are weakly assortative across protein interaction edges,²⁷
28 and only a few detected modules exhibit any disassortative signal, and even then²⁸
29 their internal r is relatively close to zero (Fig. S3). In this situation, the smooth²⁹
30 filter typically outperforms the sharp filter (Fig. 2A). ³⁰

31 We used the method proposed by Ronen and Akalin to optimize the smoothing³¹
32 parameter for the diffusion based methods by maximizing the entropy of points³²
33 embedded in a 2-dimensional PCA space [14]. Since the distributions of the healthy³³

¹tissue and delta vector data are quite different, we optimized the smoothing param-¹
²eters of each individually. ²

³ Across model complexities, we find that denoising before model training using any ³
⁴type of network filter or diffusion-based method provides a substantial reduction in ⁴
⁵prediction error relative to training on unfiltered data (Fig. 3B, Fig. S4). **The median** ⁵
⁶**filter and netSmooth have very similar performance with around 22% improvement** ⁶
⁷**in MAE from no filter. The Laplacian exponential diffusion kernel, patchwork filter** ⁷
⁸**paired with spectral community detection, and mean network filters have the lowest** ⁸
⁹**MAE, improving upon the raw data by 32%, 37%, and 43%, respectively.** ⁹

¹⁰ Error rates tend to decrease with greater model complexity K , suggesting that ¹⁰
¹¹more complex models are better able to capture variations in the precise expres- ¹¹
¹²sion level changes between tissue states. This decrease in error also occurs without ¹²
¹³first filtering the expression data. However, the improvement in prediction accuracy ¹³
¹⁴from increasing the model complexity without filtering is modest (5.2% at $K = 6$) ¹⁴
¹⁵compared to the improvement from first applying the best network filter (42% at ¹⁵
¹⁶ $K = 1$, and 43% at $K = 6$). ¹⁶

¹⁷ We note that in this real-world setting, the patchwork filter, which first partitions ¹⁷
¹⁸the protein interaction network into protein groups, **performs better with the spec-** ¹⁸
¹⁹**tral or modularity maximization algorithms than with the DC-SBM. The patchwork** ¹⁹
²⁰**filter paired with these algorithms performed very well, but the mean filter still per-** ²⁰
²¹**formed better than them. This behavior suggests that the partitions produced by** ²¹
²²**the community detection algorithms did not correlate sufficiently strongly with the** ²²
²³**underlying variation in biological signals to correctly localize the most relevant adja-** ²³
²⁴**cent measurements, in contrast to our controlled experiments** (Fig. 2D). Developing ²⁴
²⁵community detection algorithms that choose more biologically relevant partitions ²⁵
²⁶may be a useful direction of future work. ²⁶

²⁷ **Discussion** ²⁷

²⁸ Large data sets of biological signals, **such as system-wide measurements of molecular** ²⁸
²⁹**concentrations**, are often noisy. However, these measurements are not fully indepen- ²⁹
³⁰dent because they reflect the dynamics of a single interconnected system. Using a ³⁰
³¹network to represent the underlying biological relationships among a set of measure- ³¹
³²ments, we can leverage the size of these data sets to systematically denoise many ³²
³³ ³³

1 measurements at once, improving the data's utility for understanding the struc-1
2 ture and dynamics of complex biological systems or making accurate predictions in 2
3 systems biology. 3

4 Experiments using synthetic data with realistic biological network structures and 4
5 a variety of underlying signals indicates that network filters can substantially reduce 5
6 noise in large biological data sets across a broad range of circumstances (Fig. 2A,D, 6
7 Fig. S1). The greatest benefit is obtained when the type of filter is matched to the 7
8 underlying relationship among the signals, e.g., smoothing for assortative signals 8
9 (correlation) and sharpening for disassortative signals (anti-correlation). However, 9
10 for modest levels of correlation, even the wrong kind of filter yields some benefit 10
11 because of a regression to the mean effect, in which combining several related signals 11
12 filters out more noise than it introduces through bias. When signal types are het- 12
13 erogeneous across the network, so that the strength or direction of the correlation 13
14 differs in different parts of the network, a "patchwork" filter often performs bet- 14
15 ter. In this approach, we first partition the network into smaller, more homogeneous 15
16 modules (groups of interrelated measurements) and then apply filters independently 16
17 to the measurements now localized within each module (Fig. 2D). 17
18

19 In a more realistic setting, in which we train a machine learning algorithm to 19
20 predict changes in human protein expression levels when healthy tissue becomes 20
21 cancerous, applying a network filter based on a high-quality protein interaction 21
22 network before model training substantially improves prediction accuracy, compared 22
23 to training on unfiltered data (Fig. 3B). In this experiment, the protein interaction 23
24 network itself is not noise-free [29], indicating that filtering using an imperfect 24
25 network can be better than not filtering at all. Our experiment on rewiring net- 25
26 work edges further supports that network filters still work well on noisy network 26
27 structures (Fig. S1). 27

28 In each experiment, we compared our network filters to techniques relying on 28
29 network diffusion algorithms to weight the nodes before combining them. Both 29
30 netSmooth and the Laplacian exponential diffusion kernel have similar character- 30
31 istics to the smoothing network filters. In the non-modular synthetic graphs, they 31
32 perform better with more assortative underlying data. However, on modular graphs 32
33 with heterogeneous data values, the performance only slightly increases as more 33

¹communities have assortative data values, and decreases when communities have¹
²disassortative values. ²

³ We find an apparent trade-off between the size of the local area of nodes used to ³
⁴denoise a value and the range of values that can be recovered. The diffusion-based⁴
⁵techniques outperform local smoothing network filters when there is no correlation⁵
⁶or anti-correlation between neighboring values. This improvement is caused by a ⁶
⁷larger regression to the mean effect from using many more neighbors to denoise⁷
⁸any given value. While this effect is beneficial in the experiment with non-modular⁸
⁹synthetic graphs, it strongly hinders their performance on modular graphs with⁹
¹⁰heterogenous data values because the diffusion-based techniques tend to use values¹⁰
¹¹outside their community, which are drawn from different underlying distributions. ¹¹
¹²Furthermore, increasing the smoothing parameter increases the weight of values¹²
¹³outside of the community, strongly deteriorating their performance. Thus, regression¹³
¹⁴to the mean is not beneficial in this experiment since each community has a different¹⁴
¹⁵distribution of data values. On the other hand, the mean and median filters are¹⁵
¹⁶more localized and hence make fewer errors due to combining neighbors' values¹⁶
¹⁷from different communities. ¹⁷

¹⁸ There are a number of potentially valuable directions for future work on network¹⁸
¹⁹filters, which may improve their error rates or adapt them to more complicated¹⁹
²⁰settings or tasks. Techniques from image processing, both simple and advanced,²⁰
²¹represent a particularly promising direction to explore [36, 37, 38]. For instance,²¹
²²here, we only considered the network filters combine measurements associated with²²
²³directly adjacent nodes. As a result, the denoised values associated with low degree²³
²⁴nodes in the network derive from a relatively smaller number of measurements,²⁴
²⁵and hence are likely to have larger residual noise than will higher degree nodes.²⁵
²⁶Modifying the network filter for low degree nodes to look beyond nearest neighbors,²⁶
²⁷e.g., to ensure a minimum number of pooled measurements per node, may provide²⁷
²⁸better guarantees on the accuracy of the denoised value. An example of this type²⁸
²⁹of technique in image processing include the Gaussian filter [39]. ²⁹

³⁰ Image segmentation, in which an image is first partitioned into visually distinct ³⁰
³¹pieces, e.g., separating the foreground from the background, is a common prepro-³¹
³²cessing step in image analysis. The patchwork filter considered here is a simple ³²
³³adaptation of this idea, but it relies on off-the-shelf community detection algorithms ³³

to partition the nodes, considers different modules independently, and ignores connections that run between modules. While this approach should retain the most informative relationships among the measurements it also serves to reduce the degrees of many nodes, which may lessen the benefits of filtering, as described above. Furthermore, the patchwork filter will not work well on networks with disassortative community structure where nodes in the same community tend to not form edges between each other. In such cases, the patchwork filter would significantly reduce the degree of all nodes and limit the potential for network filters to denoise their data. Thus, the patchwork filter may perform best with community detection algorithms that return assortative community structures and sever the least number of edges within communities.

Developing filters that utilize the edges between modules could mitigate the induced low-degree effects that come from applying a patchwork filter to account for signal heterogeneity in the system. Such between-module edges should likely be considered separately from within-module edges, e.g., by adjusting their weights w_{ij} to more accurately capture the character of the particular signal relationship between the modules containing nodes i and j .

The benefits of a patchwork filter necessarily depends on how closely the network partition correlates with the underlying biological structure of the system. Off-the-shelf community detection algorithms may not always provide such partitions [40]. While the DC-SBM was able to recover partitions that were good for denoising in the synthetic data task, it did not perform as well as the modularity maximization and spectral algorithms on the real world data example. Since the assortativity coefficients for the Human Protein Atlas range from 0 to 0.1, the benefit is dominated by the regression to the mean effect, which does better on higher degree nodes to reduce the noise. Thus, the community detection method that finds partitions optimal for denoising may differ network to network [27]. Trying a few different community detection methods like we did here should aid in finding network partitions that best correlate with the system's underlying structure. In some settings, developing application-specific partitioning algorithms, or algorithms that can exploit biologically meaningful node attributes [26], may improve the behavior of a patchwork filter. For data sets where the data is relatively homogenous, a

1 smoothing or sharpening filter applied to the network as a whole may provide more
2 benefits than the patchwork filter.

3 Finally, the network filters defined here make few specific assumptions about the
4 underlying noise-generating process itself. In specific applications, much more may
5 be known about the direction, magnitude, and clustering of errors across large-scale
6 measurements. For instance, in molecular profiling data, endogenous biological fac-
7 tors like cell cycle effects likely induce distinct noise patterns compared to exogenous
8 technical factors like sample preparation or instrument variation. Developing more
9 application specific error models that could be combined with network filters may
10 provide more powerful denoising techniques than the general filters described here.

11

12 Conclusion

13 Network filters are a flexible tool and can exploit a variety of network data, includ-
14 ing networks of molecular binding interactions. Network filters can be extended to
15 exploit information about the sign or strength of interactions or to allow the type
16 of interaction to vary across different modules within the network. These filters
17 can also be applied to networks of any size, ranging from local signaling pathways
18 to entire protein interaction networks. In fact, any network that correlates with
19 the underlying causal structure of a set of measured variables could potentially be
20 used as a filter. By exploiting these underlying relationships, a network filter pools
21 correlated information, which mitigates independent noise, in much the same way
22 that image processing techniques use information from nearby pixels to denoise an
23 image. Overall, our study demonstrates that network filters have the potential to
24 improve the analysis of system-level biological data.

25

26 Methods

27 Synthetic data with known noise and structure

28 In the first experiment we generate simple non-modular random graphs using the
29 Chung-Lu (CL) model [41, 42, 43] with $n = 100$ nodes and a degree distribution
30 that, in expectation, follows a power law distribution $\Pr(k) \propto k^{-\alpha}$ with parameter
31 $\alpha = 3$ for $k \geq 1$. If the generated degree sequence included a node with degree
32 $k > 17$, a new degree sequence was sampled. This choice ensured that no star-
33 like subgraphs were created. In our analysis, only nodes in the largest connected

¹component were included. This choice mitigates the bias experienced by low degree¹
²nodes, which are the most likely nodes to exist outside the largest component. ²

³ For each CL synthetic network, we generate node values using the procedure³
⁴described below. We vary the assortativity coefficient $r \in [-0.8, 0.8]$ while drawing⁴
⁵values from a Normal distribution with mean and variance $\mu = \sigma^2 = 100$. We⁵
⁶simulate measurement noise by taking a random permutation of a uniformly random⁶
⁷25% of the node values. We then apply each of the networks filters (mean, median,⁷
⁸sharp) to these noisy values, and calculate the mean absolute error (MAE) of the⁸
⁹original and denoised values. We also apply netSmooth and Laplacian exponential⁹
¹⁰kernel methods varying smoothing parameter values to this data, and calculate the¹⁰
¹¹MAE of original and denoised values. Results are averaged over 5000 repetition of¹¹
¹²this process. ¹²

¹³ To create noisy non-modular networks, we performed a random rewiring procedure¹³
¹⁴previously described [30]. After generating a non-modular random graph using the¹⁴
¹⁵CL model and generating metadata, we select a given proportion of edges to remove¹⁵
¹⁶from the graph. Then, we placed the same number of new edges between any two¹⁶
¹⁷nodes chosen uniformly at random, while ensuring that there were no multi-edges¹⁷
¹⁸in the graph. Then the filters were applied to the noisy network as normal. ¹⁸

¹⁹ In the second experiment, we generate simple modular random graphs using the¹⁹
²⁰degree-corrected stochastic block model (DC-SBM) [25], with $\kappa = 5$ communities of²⁰
²¹ $n_r = 100$ nodes each ($n = 500$ nodes total), and the same degree distribution as the²¹
²²non-modular case. The network's modular structure is specified using the standard²²
²³"planted partition" model [25], in which the community mixing matrix ω_{rs} is given²³
²⁴by a linear combination of a perfectly modular graph and a random graph, and has²⁴
²⁵the form $\omega_{rs} = \lambda \omega_{rs}^{\text{planted}} + (1 - \lambda) \omega_{rs}^{\text{random}}$, with $\lambda = 0.85$. ²⁵

²⁶ For each DC-SBM network, we generate node values with the following properties: ²⁶

²⁷(i) the distribution of values within each module are drawn from a module-specific²⁷
²⁸Normal distribution with mean $\mu = \{110, 80, 60, 40, 20\}$ and variance $\sigma^2 = 25$,²⁸
²⁹(ii) $\kappa' \in [0, 5]$ communities are assigned to have negative assortativity coefficients,²⁹
³⁰and (iii) the within-community assortativity coefficients are chosen uniformly at³⁰
³¹random on the interval $|r| \in [0.4, 0.7]$. These choices construct a hard test in which³¹
³²a filter's accuracy is effectively penalized if it uses nodes outside a given community³²
³³to denoise a particular value. For the patchwork filter, we partition the network³³

¹using three different classes of community detection algorithms. The “metadata-¹
²aware” DC-SBM [26] and spectral algorithm [24] partitioned the graph in $\hat{\kappa} = 5$ ²
³communities. Modularity maximization partitioned the graph into the number of³
⁴clusters that maximizes the modularity function [23]. Noise is induced and accuracy⁴
⁵is assessed as in the non-modular case, except that the nodes are randomly permuted⁵
⁶within each module rather than the whole network. 6

⁸Generating synthetic correlated measurements 8

⁹We generate node values with a specified assortativity coefficient r_* , for a specified⁹
¹⁰adjacency matrix A , using Markov chain Monte Carlo (MCMC). The assortativity¹⁰
¹¹coefficient r is defined as 11

$$r = \frac{\sum_{ij} (A_{ij} - k_i k_j / 2m) x_i x_j}{\sum_{ij} (k_i \delta_{ij} - k_i k_j / 2m) x_i x_j}$$

¹² 12
¹³ 13
¹⁴ 14
¹⁵where k_i is the degree of node i , x_i is the value associated with node i , $2m = \sum_{ij} A_{ij}$ ¹⁵
¹⁶is twice the number of edges in the network, A_{ij} is the entry in the adjacency matrix¹⁶
¹⁷for nodes i and j , and δ_{ij} is the Kronecker delta function. 17

¹⁸ 18
¹⁹Given a network A , a desired assortativity coefficient r_* , and a node value distri-¹⁹
²⁰bution $\Pr(x)$, we generate a set of node values as follows. 20

- ²¹ 1 Assign each node a value drawn iid from $\Pr(x)$. 21
- ²² 2 Calculate the current assortativity coefficient r_0 . 22
- ²³ 3 Set $t = 1$. 23
- ²⁴ 4 While the difference between the desired and current assortativity coefficient²⁴
²⁵ $\Delta = |r_t - r_*| > \beta$, a specified tolerance, do: 25
 - ²⁶ • Pick a node i uniformly at random and assign it a new value x'_i drawn²⁶
²⁷ iid rom $\Pr(x)$. 27
 - ²⁸ • Calculate the corresponding assortativity coefficient r_t and difference²⁸
²⁹ $\Delta' = |r_t - r_*|$. 29
 - ³⁰ • If the new value does not improve the assortativity, i.e., $\Delta' > \Delta$ restore³⁰
³¹ x_i . Otherwise, increment t . 31
- ³² 5 Return the node values \mathbf{x} with the desired assortativity coefficient, r_* . 32

³³In our experiments, we set $\beta = 0.009$. 33

¹Diffusion-based denoising methods 1

²We benchmark the network filters against two comparison methods that weight 2

³nodes based on different diffusion kernels. The Laplacian exponential diffusion ker- 3

⁴nel [15] \mathbf{S}_β is defined as 4

5 5

$$6 \quad \mathbf{S}_\beta = e^{-\beta \mathbf{L}} \quad 6$$

7 7

⁸where β is a real valued smoothing parameter, and \mathbf{L} is the graph Laplacian. The 8

⁹denoised data vector is found by multiplying the matrix and noisy data vector 9

10 10

$$11 \quad \mathbf{x}_{\text{Laplacian},\beta} = \mathbf{S}_\beta \mathbf{x} \ . \quad 11$$

¹²The netSmooth method [14] uses the personalized Pagerank [28] vector to weight 12

¹³each node. This kernel \mathbf{K} is defined as 13

14 14

$$15 \quad \mathbf{K}_\alpha = (1 - \alpha)(\mathbf{I} - \alpha \mathbf{B})^{-1} \quad 15$$

16 16

¹⁷where \mathbf{B} is a adjacency matrix that is degree normalized by column such that 17

¹⁸ $\mathbf{B}_{ij} = \frac{1}{k_j}$ if there exists an edge between i and j , α is the smoothing parameter 18

¹⁹(also known as the restart probability), and \mathbf{I} is the identity matrix. The denoised 19

²⁰data vector is found by multiplying this kernel and the noisy kernel as such 20

21 21

$$22 \quad \mathbf{x}_{\text{netSmooth},\alpha} = \mathbf{K}_\alpha \mathbf{x} \ . \quad 22$$

23 23

For simplicity, we call both β and α “smoothing parameter” throughout as they 24

²⁴have a similar function for their respective methods. 24

25 25

²⁶Human protein expression and interaction 26

²⁷Protein expression data were drawn from the Human Protein Atlas (HPA) version 27

²⁸16 [33], which details protein expression in human tissues by large scale immuno- 28

²⁹histochemistry, for over 12,000 proteins in 20 tissue types, each in a healthy and 29

³⁰cancerous state. Human protein interaction (PPIN) data were drawn from the HINT 30

³¹database [35], which combines data from several interactome databases and is cu- 31

³²rated for biological interactions with high levels of evidence. The HINT network 32

³³contains $n = 12,864$ proteins and $m = 62,435$ undirected, unweighted edges. 33

¹ To construct the network filter, we first map the data from the HPA to the PPIN.¹
² HPA proteins are indexed by their Ensembl IDs, while HINT proteins are indexed²
³ by their Uniprot IDs. A map from Ensembl IDs to Uniprot IDs was constructed³
⁴ using the HGNC BioMart tool. If a node had multiple mapped expression values,⁴
⁵ we averaged them. We allow protein expression values from HPA to map to multiple⁵
⁶ nodes if the Ensembl ID maps to multiple nodes in the PPIN. If the gene expression⁶
⁷ value does not map to any nodes in the PPIN, we discard these as they cannot be⁷
⁸ de-noised by the network filters. There is one protein in the cancer dataset and 283⁸
⁹ proteins in the healthy tissue dataset missing protein expression values in no more⁹
¹⁰ than 2 cancers or healthy tissues. For these cases, we impute the missing data from¹⁰
¹¹ the same protein in another cancer or healthy tissue uniformly at random (impute¹¹
¹² healthy from healthy, and cancer from cancer).¹²

¹³ After keeping the largest connected component of nodes with associated HPA¹³
¹⁴ data values, these preprocessing steps produce a network with $n = 8,199$ proteins¹⁴
¹⁵ with IHC expression information across all 20 tissue types and both healthy and¹⁵
¹⁶ cancerous states, and $m = 37,607$ edges. The included healthy-cancerous tissue¹⁶
¹⁷ pairs are: breast, glioma, cervix, colorectal, endometrial, testicular, thyroid, renal,¹⁷
¹⁸ liver, lung, lymphoma, pancreas, prostate, skin, stomach, melanocyte, urinary, head¹⁸
¹⁹ and neck, ovary, carcinoid. For the healthy tissues, the protein expression values of¹⁹
²⁰ specific cells types that can give rise to the corresponding cancer were averaged²⁰
²¹ together to form one vector (Table S1).²¹

²² ²³ ²⁴ Predicting expression changes in human cancer ²² ²³ ²⁴

²⁵ The machine learning task is to predict the changes in protein expression levels²⁵
²⁶ when a human tissue changes types from healthy to cancerous. We use K -nearest²⁶
²⁷ neighbors regression to learn a model that can predict these changes when given²⁷
²⁸ the expression levels of a healthy tissue (Fig. 4). We train and evaluate the model²⁸
²⁹ using leave-one-out cross validation, in which the model is trained on the observed²⁹
³⁰ changes in 19 healthy-cancerous tissue pairs, and is tested on one unobserved pair.³⁰
³¹ We first train and evaluate the model on unfiltered data, and then compare its³¹
³² accuracy to a model where we apply a network filter to the expression data prior³²
³³ to training.³³

1 First, we applied principal component analysis (PCA) on the training set of 19¹
2 healthy tissue protein vectors as a feature extraction method. Next, using the em-²
3 bedded PCA space learned from the training set, we project the held-out healthy³
4 sample into the same PCA space. We then determine the K -nearest neighbors of⁴
5 the held-out healthy tissue by calculating the Euclidean distance of the first four⁵
6 principal components between this point and all other healthy tissues. ⁶

7 ⁷
8 Given this identification of which healthy tissues are most similar to the left-out ⁸
9 healthy tissue, we predict the protein expression changes for the held-out observa- ⁹
10 tion. We calculate the expression changes between cancerous and healthy tissues, ¹⁰
11 which we call a “delta” vector. Then, we perform PCA on the 19 delta vectors to ¹¹
12 extract features. The weighted average of the delta vectors corresponding to the ¹²
13 K -nearest neighbors learned from the healthy tissues are averaged together, where ¹³
14 the weight is proportional to the inverse of the Euclidean distance to the held-out ¹⁴
15 healthy tissue. Finally, we project the predicted delta vector from four principal ¹⁵
16 components back to the $n = 8,199$ proteins and calculate the mean absolute error ¹⁶
17 (MAE) of this vector and the actual delta vector. ¹⁷

18 The basic networks filters evaluated in this task have the form given in the ¹⁸
19 main text. For the patchwork filter, the DC-SBM or spectral algorithms partition ¹⁹
20 the PPIN into $\kappa = 10$ communities, and modularity maximization automatically ²⁰
21 chooses the number of communities that maximizes the modularity function. Then, ²¹
22 we apply the mean filter within each community. ²²

23 For the diffusion-based methods, we choose optimized smoothing parameters for ²³
24 the human protein expression data set using a method described by the netSmooth ²⁴
25 authors [14] to maximize the entropy of a 2D embedding of the data. As the ²⁵
26 healthy data and delta vectors have different data distributions, we choose opti- ²⁶
27 mized smoothing parameters for each data set separately. Briefly, the healthy tissue ²⁷
28 protein expression or delta vectors were embedded in a PCA space with the first ²⁸
29 two principal components. This space was discreteized into a four by four grid, ²⁹
30 equally spaced from the data points at the minimum and maximum of each PC. We ³⁰
31 calculated the Shannon entropy, $H(x) = -\sum_i P(x_i) \log P(x_i)$, of this discretized ³¹
32 embedding, and chose the smallest smoothing parameter that maximized the en- ³²
33 tropy. For netSmooth, the smoothing parameter was 0.2 for the healthy tissues and ³³

¹ 0.3 for the delta vectors. And for the Laplacian exponential diffusion kernel, it was ¹	
² 0.2 and 0.1 for healthy tissues and delta vectors, respectively.	2
³	3
⁴Abbreviations	4
⁵ CL: Chung-Lu	5
⁶ DC-SBM: Degree-Corrected Stochastic Block Model	6
⁷ GSEA: Gene Set Enrichment Analysis	7
⁸ HGNC: HUGO Gene Nomenclature Committee	8
⁹ HINT: High-Quality Interactomes	9
¹⁰ HPA: Human Protein Atlas	10
¹¹ IHC: Immunohistochemistry	11
¹² MAE: Mean Absolute Error	12
¹³ MCMC: Markov Chain Monte Carlo	13
¹⁴ PCA: Principal Component Analysis	14
¹⁵ PPIN: Protein-Protein Interaction Network	15
¹⁶	16
¹⁷Declarations	17
Ethics approval and consent to participate	
¹⁸ Not applicable	18
¹⁹Consent for publication	19
²⁰ Not applicable	20
²¹Availability of data	21
The datasets and code supporting the conclusions of the article are available in github repository	
²² https://github.com/andykavran/Network_Filters	22
²³Competing interests	23
²⁴ The authors declare that they have no competing interests.	24
²⁵Funding	25
²⁶ AJK was supported in part by the Interdisciplinary Quantitative Biology (IQ Biology) Program (NSF IGERT grant number 1144807) at the BioFrontiers Institute, University of Colorado, Boulder. AC was supported in part by Grant	26
²⁷ No. IIS-1452718 from the National Science Foundation. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.	27
²⁸Author's contributions	28
²⁹ Conceptualization: AC, AJK	29
Data Curation: AC, AJK	
³⁰ Formal Analysis: AJK	30
³¹ Funding Acquisition: AC	31
Investigation: AC, AJK	
³² Software: AJK	32
Visualization: AJK	
³³ Writing - original draft: AC, AJK	33

¹ Writing - review and editing: AC, AJK	1
² Both authors have read and approved the manuscript.	2
³ Acknowledgements	3
We would like to thank Natalie Ahn and Mark Newman for helpful conversations.	
⁴	4
Author details	
⁵ ¹ Department of Biochemistry, University of Colorado, Boulder, CO, USA. ² BioFrontiers Institute, University of	5
⁶ Colorado, Boulder, CO, USA. ³ Department of Computer Science, University of Colorado, Boulder, CO, USA.	6
⁴ Santa Fe Institute, Santa Fe, NM, USA.	
⁷	7
References	
⁸ 1. Woodworth, M.B., Girsakis, K.M., Walsh, C.A.: Building a lineage from single cells: genetic techniques for cell	8
lineage tracking. <i>Nature Reviews Genetics</i> 18 (4), 230 (2017)	9
⁹ 2. McKenna, A., Gagnon, J.A.: Recording development with single cell dynamic lineage tracing. <i>Development</i>	9
¹⁰ 146 (12), 169730 (2019)	10
¹¹ 3. Pastushenko, I., Blanpain, C.: Emt transition states during tumor progression and metastasis. <i>Trends in Cell</i>	11
<i>Biology</i> (2018)	11
¹² 4. Hugo, W., Shi, H., Sun, L., Piva, M., Song, C., Kong, X., Moriceau, G., Hong, A., Dahlman, K.B., Johnson,	12
¹³ D.B., <i>et al.</i> : Non-genomic and immune evolution of melanoma acquiring mapki resistance. <i>Cell</i> 162 (6),	13
1271–1285 (2015)	13
¹⁴ 5. Muranen, T., Selfors, L.M., Worster, D.T., Iwanicki, M.P., Song, L., Morales, F.C., Gao, S., Mills, G.B.,	14
¹⁵ Brugge, J.S.: Inhibition of pi3k/mtor leads to adaptive resistance in matrix-attached cancer cells. <i>Cancer Cell</i>	15
21 (2), 227–239 (2012)	15
¹⁶ 6. Michiels, S., Koscielny, S., Hill, C.: Prediction of cancer outcome with microarrays: a multiple random	16
validation strategy. <i>The Lancet</i> 365 (9458), 488–492 (2005)	16
¹⁷ 7. Button, K.S., Ioannidis, J.P., Mokrysz, C., Nosek, B.A., Flint, J., Robinson, E.S., Munafò, M.R.: Power failure:	17
¹⁸ why small sample size undermines the reliability of neuroscience. <i>Nature Reviews Neuroscience</i> 14 (5), 365	18
(2013)	18
¹⁹ 8. Mootha, V.K., Lindgren, C.M., Eriksson, K.F., Subramanian, A., Sihag, S., Lehar, J., Puigserver, P., Carlsson,	19
E., Ridderstrale, M., Laurila, E., Houstis, N., Daly, M.J., Patterson, N., Mesirov, J.P., Golub, T.R., Tamayo, P.,	20
²⁰ Spiegelman, B., Lander, E.S., Hirschhorn, J.N., Altshuler, D., Groop, L.C.: PGC-1 α -responsive genes involved	20
²¹ in oxidative phosphorylation are coordinately downregulated in human diabetes. <i>Nature Genetics</i> 34 (3),	21
267–273 (2003)	21
²² 9. Barbie, D.A., Tamayo, P., Boehm, J.S., Kim, S.Y., Moody, S.E., Dunn, I.F., Schinzel, A.C., Sandy, P., Meylan,	22
²³ E., Scholl, C., <i>et al.</i> : Systematic rna interference reveals that oncogenic kras-driven cancers require tbk1.	23
<i>Nature</i> 462 (7269), 108 (2009)	23
²⁴ 10. Ronan, T., Qi, Z., Naegle, K.M.: Avoiding common pitfalls when clustering biological data. <i>Science Signaling</i>	24
9 (432), 6–6 (2016)	24
²⁵ 11. Sørlie, T., Perou, C.M., Tibshirani, R., Aas, T., Geisler, S., Johnsen, H., Hastie, T., Eisen, M.B., Van De Rijn,	25
²⁶ M., Jeffrey, S.S., <i>et al.</i> : Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical	26
implications. <i>Proceedings of the National Academy of Sciences USA</i> 98 (19), 10869–10874 (2001)	26
²⁷ 12. Gupta, A., Wang, H., Ganapathiraju, M.: Learning structure in gene expression data using deep architectures,	27
²⁸ with an application to gene clustering. In: <i>Bioinformatics and Biomedicine (BIBM)</i> , 2015 IEEE International	28
Conference On, pp. 1328–1335 (2015). IEEE	28
²⁹ 13. Tan, J., Ung, M., Cheng, C., Greene, C.S.: Unsupervised feature construction and knowledge extraction from	29
³⁰ genome-wide assays of breast cancer with denoising autoencoders. In: <i>Pacific Symposium on Biocomputing</i> , pp.	30
132–143 (2015)	30
³¹ 14. Ronen, J., Akalin, A.: netsmooth: Network-smoothing based imputation for single cell rna-seq. <i>F1000Research</i>	31
7 (2018)	31
³² 15. Dørnum, G., Snipen, L., Solheim, M., Saebø, S.: Smoothing gene expression data with network information	32
³³ improves consistency of regulated genes. <i>Statistical applications in genetics and molecular biology</i> 10 (1) (2011)	33
16. Newman, M.E.J.: Mixing patterns in networks. <i>Physical Review E</i> 67 (2), 026126 (2003)	33

17. Ideker, T., Ozier, O., Schwikowski, B., Siegel, A.F.: Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics* **18**(suppl_1), 233–240 (2002)
18. Goncalves, A., Leigh-Brown, S., Thybert, D., Stefflova, K., Turro, E., Flicek, P., Brazma, A., Odom, D.T., Marioni, J.C.: Extensive compensatory cis-trans regulation in the evolution of mouse gene expression. *Genome Research* **22**(12), 2376–2384 (2012)
19. Bauer, P.M., Fulton, D., Bo, Y.C., Sorescu, G.P., Kemp, B.E., Jo, H., Sessa, W.C.: Compensatory phosphorylation and protein-protein interactions revealed by loss of function and gain of function mutants of multiple serine phosphorylation sites in endothelial nitric-oxide synthase. *The Journal of Biological Chemistry* **278**(17), 14841–14849 (2003)
20. Peel, L., Delvenne, J.-C., Lambiotte, R.: Multiscale mixing patterns in networks. *Proceedings of the National Academy of Sciences USA* **115**(16), 4057–4062 (2018)
21. Rudolph, J.D., de Graauw, M., van de Water, B., Geiger, T., Sharan, R.: Elucidation of signaling pathways from large-scale phosphoproteomic data using protein interaction networks. *Cell Systems* **3**(6), 585–593 (2016)
22. Mansourpour, M., Rajabi, M., Blais, J.: Effects and performance of speckle noise reduction filters on active radar and sar images. In: *Proc. ISPRS*, vol. 36, p. 41 (2006)
23. Blondel, V.D., Guillaume, J.-L., Lambiotte, R., Lefebvre, E.: Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment* **2008**(10), 10008 (2008)
24. Ng, A., Jordan, M., Weiss, Y.: On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems* **14**, 849–856 (2001)
25. Karrer, B., Newman, M.E.J.: Stochastic blockmodels and community structure in networks. *Physical Review E* **83**(1), 016107 (2011)
26. Newman, M.E.J., Clauset, A.: Structure and inference in annotated networks. *Nature Communications* **7**, 11863 (2016)
27. Ghasemian, A., Hosseinmardi, H., Clauset, A.: Evaluating overfit and underfit in models of network community structure. *IEEE Transactions on Knowledge and Data Engineering* **32**(9), 1722–1735 (2019)
28. Jeh, G., Widom, J.: Scaling personalized web search. In: *Proceedings of the 12th International Conference on World Wide Web*, pp. 271–279 (2003). *Acm*
29. Hart, G.T., Ramani, A.K., Marcotte, E.M.: How complete are current yeast and human protein-interaction networks? *Genome biology* **7**(11), 120 (2006)
30. Middendorf, M., Ziv, E., Wiggins, C.H.: Inferring network mechanisms: The drosophila melanogaster protein interaction network **102**(9), 3192–3197 (2005)
31. Campbell, J.D., Mazzilli, S.A., Reid, M.E., Dhillon, S.S., Platero, S., Beane, J., Spira, A.E.: The case for a pre-cancer genome atlas (pcga). *Cancer Prevention Research* **9**(2), 119–124 (2016)
32. Spira, A., Yurgelun, M.B., Alexandrov, L., Rao, A., Bejar, R., Polyak, K., Giannakis, M., Shilatifard, A., Finn, O.J., Dhodapkar, M., et al.: Precancer atlas to drive precision prevention trials. *Cancer Research* **77**(7), 1510–1541 (2017)
33. Uhlén, M., Fagerberg, L., Hallström, B.M., Lindskog, C., Oksvold, P., Mardinoglu, A., Sivertsson, Å., Kampf, C., Sjöstedt, E., Asplund, A., et al.: Tissue-based map of the human proteome. *Science* **347**(6220), 1260419 (2015)
34. Vyberg, M., Nielsen, S.: Proficiency testing in immunohistochemistry—experiences from nordic immunohistochemical quality control (nordiqc). *Virchows Archiv* **468**(1), 19–29 (2016)
35. Das, J., Yu, H.: Hint: High-quality protein interactomes and their applications in understanding human disease. *BMC Systems Biology* **6**(1), 92 (2012)
36. Motwani, M.C., Gadiya, M.C., Motwani, R.C., Harris, F.C.: Survey of image denoising techniques. In: *Proceedings of GSPX*, pp. 27–30 (2004)
37. Agostinelli, F., Anderson, M.R., Lee, H.: Adaptive multi-column deep neural networks with application to robust image denoising. In: *Advances in Neural Information Processing Systems*, pp. 1493–1501 (2013)
38. Öktem, R., Egiazarian, K., Lukin, V.V., Ponomarenko, N.N., Tsybal, O.V.: Locally adaptive dct filtering for signal-dependent noise removal. *EURASIP Journal on Advances in Signal Processing* **2007**(1), 042472 (2007)
39. Deng, G., Cahill, L.: An adaptive gaussian filter for noise reduction and edge detection. In: *1993 IEEE Conference Record Nuclear Science Symposium and Medical Imaging Conference*, pp. 1615–1619 (1993). *IEEE*

- 1 40. Peel, L., Larremore, D.B., Clauset, A.: The ground truth about metadata and community detection in 1
 2 networks. *Science Advances* **3**(5), 1602548 (2017) 2
 3 41. Aiello, W., Chung, F., Lu, L.: A random graph model for power law graphs. *Experimental Mathematics* **10**(1), 3
 4 53–66 (2001)
 4 42. Chung, F., Lu, L.: Connected components in random graphs with given expected degree sequences. *Annals of* 4
 5 *Combinatorics* **6**, 125–145 (2002)
 5 43. Alam, M., Khan, M., Vullikanti, A., Marathe, M.: An efficient and scalable algorithmic method for generating 5
 6 large-scale random graphs. In: *SC'16: Proceedings of the International Conference for High Performance* 6
 7 *Computing, Networking, Storage and Analysis*, pp. 372–383 (2016). IEEE 7

7 Figures 7

8 **Figure 1 Schematics of Network Filters.** Network filters are tools that denoise real-valued 8
 9 biological data using a biologically meaningful network to exploit the correlation (“smoothing”) or 9
 10 anti-correlation (“sharpening”) among neighboring measurements. **A.** A measurement x_i and its 10
 11 neighboring values in network, where the color intensity is proportional to the measured value. In 11
 12 applying the smooth filter, x_i is adjusted to be more similar to its neighbors; in applying the sharp 12
 13 filter, x_i is adjusted to be more distant from its neighbors. **B.** Measurements can also first be 13
 14 partitioned into groups (dashed line) by detecting structural modules within the network, and then 14
 15 different filters applied to different modules, ignoring between-module edges, e.g., if the signals 15
 16 are assortative in some communities and disassortative in others. 16

17 **Figure 2 Filter performance on synthetic networks.** Network filter tests on synthetic graphs with 17
 18 varying structures and known noise. The Mean Absolute Error (MAE) of **A.** network filters, **B.** 18
 19 Laplacian exponential diffusion kernel, and **C.** netSmooth on the permuted nodes as a function of 19
 20 the assortativity coefficient of 5,000 instances of noisy non-modular graphs. The smooth filters 20
 21 (mean and median) perform best on assortative data ($r > 0$), while the sharp filter is optimal for 21
 22 disassortative data ($r < 0$). When data are neither assortative nor disassortative ($r \approx 0$), 22
 23 netSmooth and Laplacian exponential kernels perform best. The MAE of **D.** network filters, **E.** 23
 24 Laplacian exponential diffusion kernel, and **F.** netSmooth on the permuted nodes as a function of 24
 25 the fraction of communities with assortative data values for 100 instances of noisy modular 25
 26 graphs. Each network instance has 5 communities and we vary how many communities have 26
 27 assortative vs. disassortative data values with a moderate assortativity coefficient $|r| \in [0.4, 0.7]$. 27
 28 The shaded areas indicate 99% bootstrapped confidence intervals. 28
 29 29
 30 30
 31 31
 32 32
 33 33

25 Additional Files 25

26 Additional file 1 — Supplemental Figures and Tables 26

27 **Figure S1 - Filter performance on rewired synthetic networks.** Figure S2 - Filter performance on modular synthetic 27
 28 networks, including the sharp filter. Figure S3 - Distribution of assortativity coefficients of network modules with 28
 29 Human Protein Atlas Data. Figure S4 - KNN regression of Human Protein Atlas data with all network filters. Table 29
 30 S1. Cell types from the Human Protein Atlas dataset averaged together to form a single healthy tissue vector. 30
 31 31
 32 32
 33 33

1
2
3
4
5 **Figure 3 Denoising to Predict Protein Expression Changes in Healthy and Cancerous Tissues.**
6 Tests of the network filters on a cancer protein expression prediction task. In this test, we predict
7 the protein expression changes that occur when a healthy tissue becomes cancerous, quantified by
8 the out-of-sample prediction accuracy with and without using network filters to preprocess the
9 data before training. **A.** The first two principal components of immunohistochemistry data of
10 healthy and cancerous tissues in the Human Protein Atlas. Arrows connect a healthy tissue (blue)
11 to the corresponding cancer (red). The first component captures variations across tissues, while
12 the second captures variation in state (healthy vs. cancerous). Predicting the precise changes
13 between healthy and cancerous tissues is a non-trivial task. Therefore, we perform a K-Nearest
14 Neighbors regression on the HPA data, with and without preprocessing with network filters. We
15 evaluate the model by leave-one-out cross validation, and calculating the MAE of the predicted
16 and actual data values for the left out healthy-cancerous pair. **B.** All network filters and diffusion
17 methods improve the MAE compared to the no-filter baseline. We compare this across different
18 choices of K , as it is a free parameter. The shaded areas represent 95% bootstrapped confidence
19 intervals

20
21
22
23 **Figure 4 Schematic of K-Nearest Neighbors Regression Framework.** We designed a weighted
24 K-nearest neighbors regression framework to predict the protein expression changes a healthy
25 tissue would undergo when becoming cancerous, given a vector of protein expression profile of a
26 healthy tissue. First, we extract features from the training set of 19 healthy tissue protein
27 expression vectors by PCA. Second, we project the left out healthy vector down to the same PCA
28 space, and third, determine K-nearest neighbors to use for the prediction task. Fourth, we extract
29 the features from the 19 delta vectors by PCA, and fifth, predict the delta vector for the left-out
30 healthy sample by taking the weighted average of the K-nearest neighbor's delta vectors. Finally,
31 sixth, we project the predicted delta vector from PCA space back to a vector of protein expression
32 values to calculate the error.

33

Figures

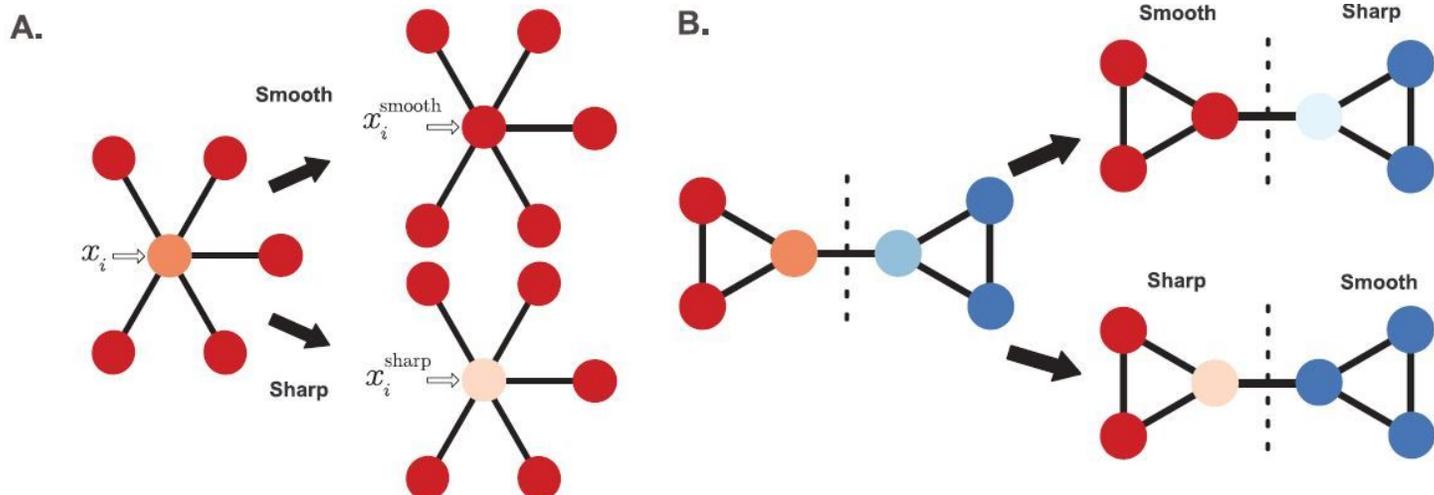


Figure 1

Schematics of Network Filters. Network filters are tools that denoise real-valued biological data using a biologically meaningful network to exploit the correlation (“smoothing”) or anti-correlation (“sharpening”) among neighboring measurements. A. A measurement x_i and its neighboring values in network, where the color intensity is proportional to the measured value. In applying the smooth filter, x_i is adjusted to be more similar to its neighbors; in applying the sharp filter, x_i is adjusted to be more distant from its neighbors. B. Measurements can also first be partitioned into groups (dashed line) by detecting structural modules within the network, and then different filters applied to different modules, ignoring between-module edges, e.g., if the signals are assortative in some communities and disassortative in others.

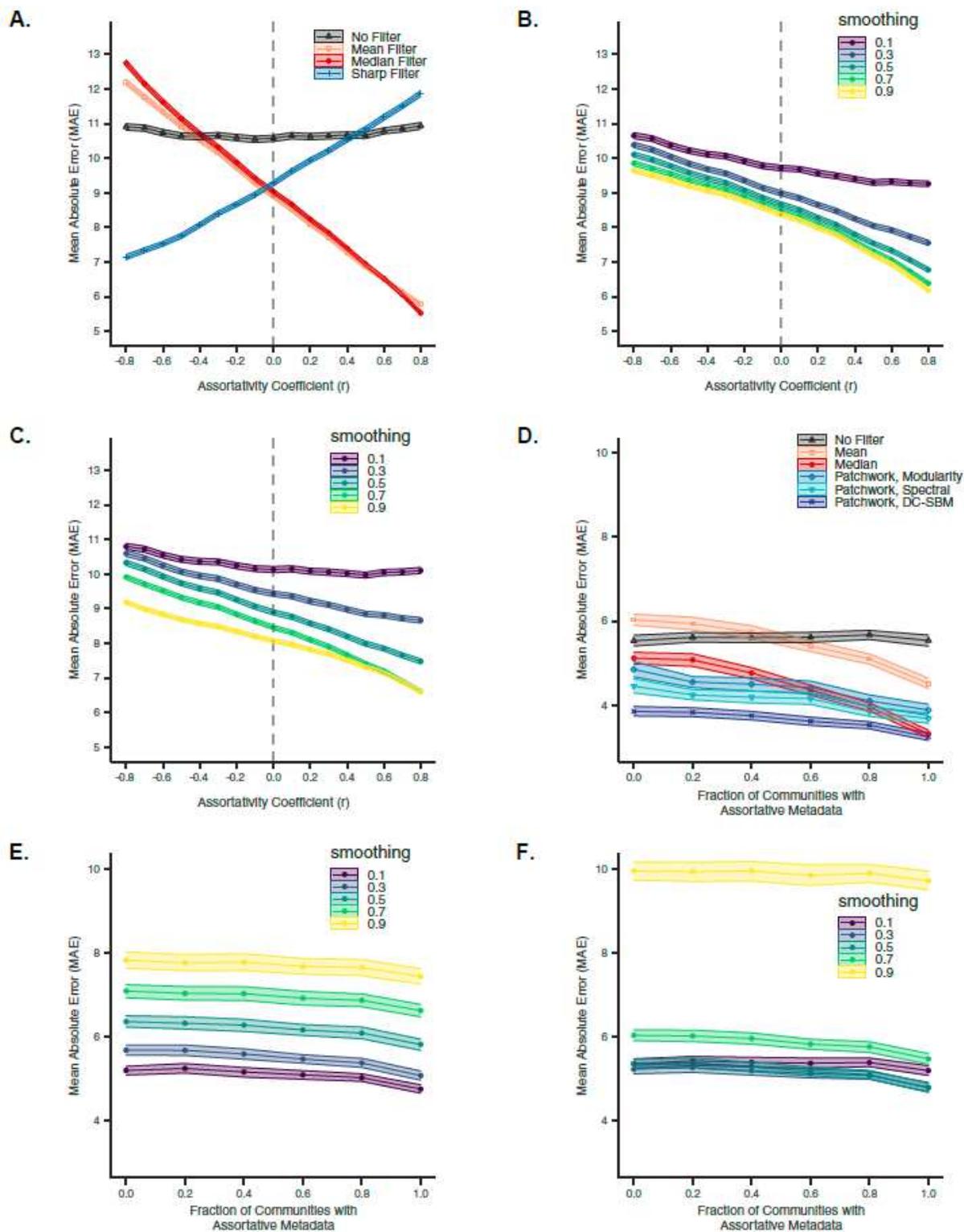


Figure 2

Filter performance on synthetic networks. Network filter tests on synthetic graphs with varying structures and known noise. The Mean Absolute Error (MAE) of A. network filters, B. Laplacian exponential diffusion kernel, and C. netSmooth on the permuted nodes as a function of the assortativity coefficient of 5,000 instances of noisy non-modular graphs. The smooth filters (mean and median) perform best on assortative data ($r > 0$), while the sharp filter is optimal for disassortative data ($r < 0$). When data are

neither assortative nor disassortative ($r \approx 0$), netSmooth and Laplacian exponential kernels perform best. The MAE of D. network filters, E. Laplacian exponential diffusion kernel, and F. netSmooth on the permuted nodes as a function of the fraction of communities with assortative data values for 100 instances of noisy modular graphs. Each network instance has 5 communities and we vary how many communities have assortative vs. disassortative data values with a moderate assortativity coefficient $|r| \in [0.4, 0.7]$. The shaded areas indicate 99% bootstrapped confidence intervals.

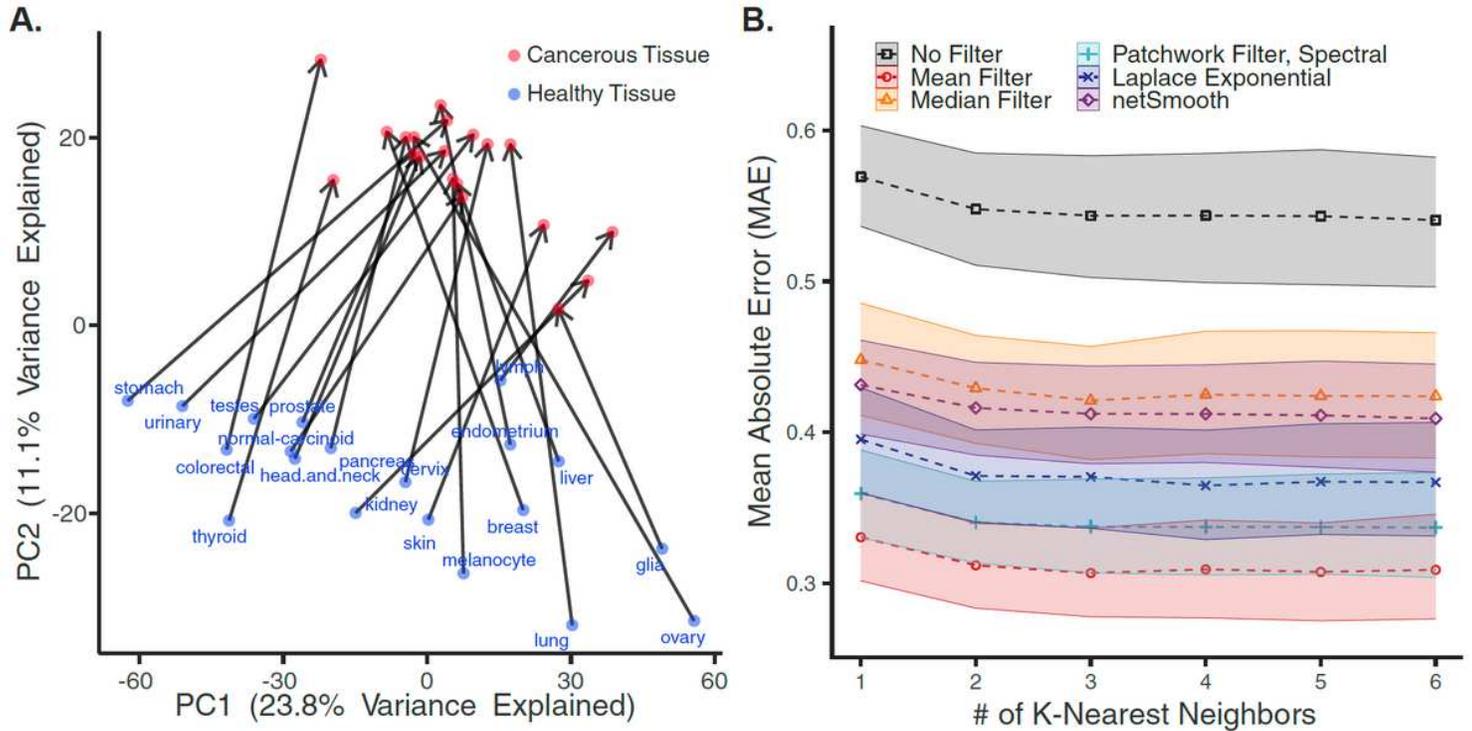


Figure 3

Denosing to Predict Protein Expression Changes in Healthy and Cancerous Tissues. Tests of the network filters on a cancer protein expression prediction task. In this test, we predict the protein expression changes that occur when a healthy tissue becomes cancerous, quantified by the out-of-sample prediction accuracy with and without using network filters to preprocess the data before training. A. The first two principal components of immunohistochemistry data of healthy and cancerous tissues in the Human Protein Atlas. Arrows connect a healthy tissue (blue) to the corresponding cancer (red). The first component captures variations across tissues, while the second captures variation in state (healthy vs. cancerous). Predicting the precise changes between healthy and cancerous tissues is a non-trivial task. Therefore, we perform a K-Nearest Neighbors regression on the HPA data, with and without preprocessing with network filters. We evaluate the model by leave-one-out cross validation, and calculating the MAE of the predicted and actual data values for the left out healthy-cancerous pair. B. All network filters and diffusion methods improve the MAE compared to the no-filter baseline. We compare this across different choices of K, as it is a free parameter. The shaded areas represent 95% bootstrapped confidence intervals

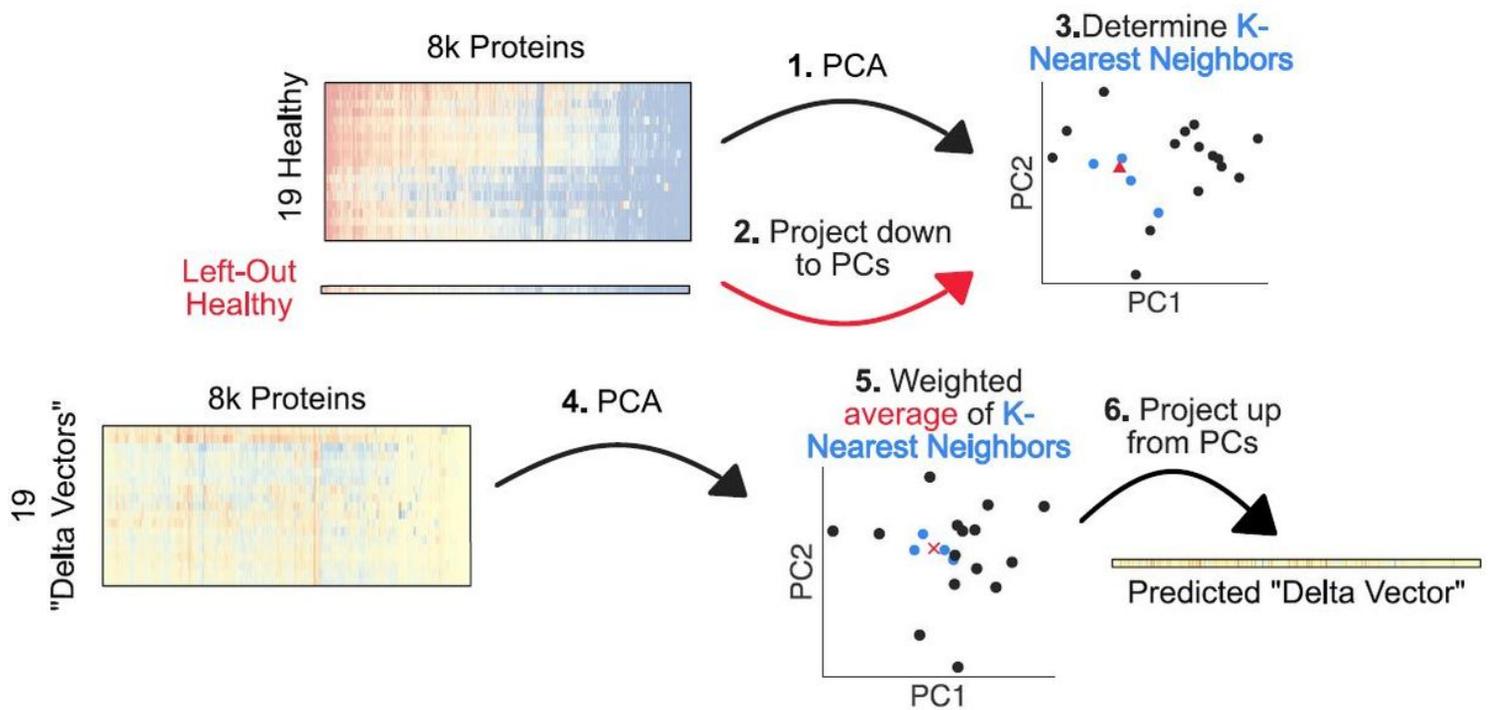


Figure 4

Schematic of K-Nearest Neighbors Regression Framework. We designed a weighted K-nearest neighbors regression framework to predict the protein expression changes a healthy tissue would undergo when becoming cancerous, given a vector of protein expression profile of a healthy tissue. First, we extract features from the training set of 19 healthy tissue protein expression vectors by PCA. Second, we project the left out healthy vector down to the same PCA space, and third, determine K-nearest neighbors to use for the prediction task. Fourth, we extract the features from the 19 delta vectors by PCA, and fifth, predict the delta vector for the left-out healthy sample by taking the weighted average of the K-nearest neighbor's delta vectors. Finally, sixth, we project the predicted delta vector from PCA space back to a vector of protein expression values to calculate the error.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [kavranclausetadditionalfile1.pdf](#)