

PNAS

Monica Nicolau^a, Arnold J. Levine^{b,1}, and Gunnar Carlsson^{a,c}

^aDepartment of Mathematics, Stanford University, Stanford, CA 94305; ^bSchool of Natural Sciences, Institute for Advanced Study, Princeton, NJ 08540; and ^cAyasdi, Inc., Palo Alto, CA 94301

Contributed by Arnold J. Levine, February 25, 2011 (sent for review July 23, 2010)

High-throughput biological data, whether generated as sequencing, transcriptional microarrays, proteomic, or other means, continues to require analytic methods that address its high dimensional aspects. Because the computational part of data analysis ultimately identifies shape characteristics in the organization of data sets, the mathematics of shape recognition in high dimensions continues to be a crucial part of data analysis. This article introduces a method that extracts information from high-throughput microarray data and, by using topology, provides greater depth of information than current analytic techniques. The method, termed *Progression Analysis of Disease (PAD)*, first identifies robust aspects of cluster analysis, then goes deeper to find a multitude of biologically meaningful shape characteristics in these data. Additionally, because *PAD* incorporates a visualization tool, it provides a simple picture or graph that can be used to further explore these data. Although *PAD* can be applied to a wide range of high-throughput data types, it is used here as an example to analyze breast cancer transcriptional data. This identified a unique subgroup of *Estrogen Receptor-positive (ER⁺)* breast cancers that express high levels of *c-MYB* and low levels of innate inflammatory genes. These patients exhibit 100% survival and no metastasis. No supervised step beyond distinction between tumor and healthy patients was used to identify this subtype. The group has a clear and distinct, statistically significant molecular signature, it highlights coherent biology but is invisible to cluster methods, and does not fit into the accepted classification of *Luminal A/B*, *Normal-like* subtypes of *ER⁺* breast cancers. We denote the group as *c-MYB⁺* breast cancer.

applied topology | p53 | systems biology

Increasingly it has become clear that, for most cancers, understanding the disease demands exploring biological processes as complex functioning systems and the pathology observed as a disruption in the coordinated performance of such systems. This viewpoint necessitates incorporating high-throughput data in the study of these diseases and consequently demands the continued development of mathematical analytic methods geared specifically to such data. The fundamental mathematical challenges in extracting meaningful information from high-throughput biological data stem, ultimately, from the difficulty in understanding the intrinsic shape of data in high dimensions (1). Shape characteristics such as kurtosis, modality, or the presence of outliers have always played a crucial role in the analysis of data, but the high dimensionality of genomic data poses mathematical difficulties in identifying its geometry. Additionally, biological phenomena are intrinsically highly variable and stochastic in nature, and notions of biological similarity are less rigid. Consequently, analysis methods for biomedical data need to identify shape characteristics that are fairly robust to changes by rescaling of distances and therefore become more qualitative in nature. This has led us to use methods adapted from the mathematics area of topology, which studies precisely the characteristics of shapes that are not rigid. The particular method we introduce in the present

article is intermediate between clustering and more distance-sensitive methods like *Principal Component Analysis (PCA)* and multidimensional scaling. This hybrid approach is able to extract unique biology from data sets. As an example, we applied our method of analysis to breast cancer transcriptional genomic data and identified a molecularly distinct unique breast cancer subgroup of *Estrogen Receptor*-positive (ER^+) tumors that have 100% overall survival and whose molecular signature is distinct from normal tissue and other breast cancers.

This article introduces *Progression Analysis of Disease (PAD)*, an approach to data analysis of disease that unravels the geometry of data sets and provides an easily accessible picture of the outcome. This method is an application of *Mapper* (2), a mathematical tool that builds a simple geometric representation of data along preassigned guiding functions called filters. *Mapper* provides both a method for mathematical data analysis and a visualization tool; the filter functions introduced through *Mapper* define a framework for supervised analysis. The output of the analysis approximates a collapse of the data into a simple, low dimensional shape, and the filter functions act as guides along which the collapse is done. *Mapper* has already been used successfully to uncover unique subtle aspects of the folding patterns of RNA (3). Here we define an application of *Mapper* to the analysis of transcriptionally genomic data from disease, with guiding filter functions provided by *Disease-Specific Genomic Analysis (DSGA)* (4). *DSGA* is a method of mathematical analysis of genomic data that highlights the component of data relevant to disease, by defining a transformation that measures the extent to which diseased tissue deviates from healthy tissue. *DSGA* has been shown to both (i) outperform traditional methods of analysis, and (ii) highlight unique biology. In combination with *Mapper*, *DSGA* transformations provide a means to define the guiding filter function, essentially by unraveling the data according to the extent of overall deviation from a healthy state.

We make *PAD* available as a Web tool, with options for *DSGA* only, *Mapper* only, or a combination of the two (5).

Our method, *PAD*, is able to identify geometric characteristics of these data that are obscured when using cluster analysis. Long gradual drifts in the graphs of these data are visible, as for example are expected when the results consist of patients with progressively advanced stages of disease. More importantly, by preserving the geometry of these data, *PAD* has identified a unique subset of breast cancers that exhibit clear and coherent clinical characteristics. Specifically, we applied *PAD* to breast cancer transcriptional microarray data (6) and identified two

Author contributions: M.N., A.J.L., and G.C. designed research; M.N. performed research; M.N., A.J.L., and G.C. analyzed data; and M.N., A.J.L., and G.C. wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. E-mail: alevine@ias.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1102826108/-/DCSupplemental.

distinct ER^+ molecular subtypes with 100% overall survival, whose molecular signatures are distinct from one another. It is important to note that survival information, given above, was not incorporated into the original analysis; rather, these two groups of patients were identified solely on the basis of gene expression data and its geometry in space. When the survival characteristics of each group were explored after *PAD* analysis was completed, each group turned out to have 100% overall survival. Both groups are ER^+ and *her2*-amplification negative ($her2^-$). One of these groups has a molecular signature that is similar to that of normal tissue and has been observed before and denoted as *Normal-like* (7). The other group is previously uncharacterized: it is composed of tumors that (i) are ER^+ and $her2^-$, (ii) express high levels of the *c-MYB* gene, (iii) express very low levels of a number of innate immune inflammatory genes, (iv) have a molecular signature that is distinct from normal tissue, and (v) do not fit into the previously accepted molecular subtypes of breast cancer (7). We have named this group the *c-MYB*⁺ group, and it constitutes 10% of *ER* tumors. This *c-MYB*⁺ group was identified and validated in an independent breast cancer data set (8).

1. Preliminary Mathematical Tools

The method consists in applying *Mapper* to genomic data from a disease state, along with the data transformation defined by *DSGA*. *Mapper* is one tool developed under the heading of topological data analysis, a recently developed form of data analysis that has a greater degree of robustness to noise and to changes in notions of distance and similarity than more distance-rigid methods like *PCA* and multidimensional scaling. Specifically, *Mapper* has the following properties: (i) its output is a combinatorial graph, rather than a linear subspace or a scattered set of points in a low-dimensional Euclidean space; (ii) the output has a multiresolution form (i.e., the data may be viewed at various scales of resolution), which is useful in distinguishing between real features and artifacts; (iii) the method has the ability to capture detail even in a large data set, in situations in which standard methods would tend to wash out the detail in question; and (iv) the method can be applied to any situation in which there is a notion of similarity or nearness, not only in Euclidean data.

1.1. Mapper. *Mapper* (2) is a mathematical tool that uses recent developments in the area of applied topology to identify shape characteristics of data sets. Topological approaches generally preserve a notion of nearness between points but can distort large-scale distances. This can be highly desirable when working with certain types of data in which, whereas small distances between points carry a notion of similarity or nearness, large distances often carry little meaning. This property often fits biological data especially well. The key idea is to identify local clusters within the data and then to understand the interaction between these small clusters by connecting them to form a graph whose shape captures aspects of the topology of the data set. *Mapper* is a mathematical tool that identifies the shape of a data set along a preassigned filter function. In its simplest form, the method works essentially as follows: we begin with a function f defined on the data and fragment the range of f into overlapping pieces. We then cluster separately the portion of the data that is mapped to each single piece. Each such local cluster can be viewed as a bin of data points. Once all data points have been assigned to bins, edges connecting bins are added: two bins that have data points in common are connected by an edge, thereby creating a graph whose shape captures important aspects of the data shape. Bins are then colored by the average value of the filter function defined on the data points inside the bin. Numeric values of these means are translated into colors, just as numeric entries in a data matrix are turned into color to produce heat maps.

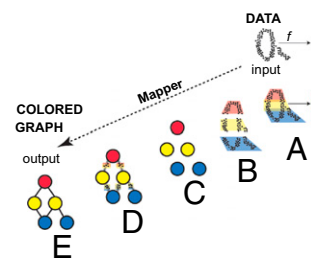


Fig. 1. *Mapper* starts with a set of data points and a filter function f and produces a colored graph that captures the shape of the data. (A) The image of the function f is subdivided into overlapping intervals. (B) Each piece is clustered separately. (C) Each cluster is represented by a colored disk: a bin of points. The color of each bin corresponds to the average value of the filter function f on the data points inside the bin. (D) Identify pairs of bins that have points in common and (E) connect pairs of bins that have points in common by an edge.

Fig. 1 illustrates how the *Mapper* construction turns a set of points with a roughly circular shape into a graph capturing this shape. *Mapper* extends a concept from topology called the nerve of a covering to the more difficult setting of working with discrete sets of points. Clearly similar shapes have similar graphs, even when the shape is somewhat distorted. However, different shapes produce different graphs that cannot be mapped into each other. Thus, *Mapper* graphs associated to data sets preserve a wealth of information about the original shapes, while providing a simplified mathematical object. Applying *Mapper* to genomic data can produce an equally simple graph from a shape that is much less accessible, because the data are both extremely high dimensional and very sparse.

1.2. Disease-Specific Genomic Analysis. *DSGA* (4) is a mathematical method for transforming omic data from diseased tissue as a sum of two terms: the *normal component* of these data best mimics healthy tissue, whereas the *disease component* measures the error or deviation from normal:

$$\vec{T} = Nc.\vec{T} + Dc.\vec{T}. \quad [1]$$

This decomposition is defined by computing a linear model of the diseased tissue data onto a *Healthy State Model (HSM)* estimated from normal tissue data, to obtain the normal component. The disease component is then the vector of error terms from the linear model fit. The *HSM* is constructed from the normal tissue data using the *FLAT* construction: a combination of mathematical data desparsing—a method to make data in very high dimensions less sparse—followed by dimension reduction through *PCA*. The *FLAT* construction was introduced by Nicolau et al. (4), and details are found in the Math Supplement of that article. Fig. 2 shows a schematic of the *DSGA* decomposition into disease and normal components. By working with the dis-

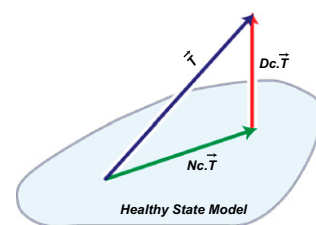


Fig. 2. *DSGA* decomposition of the original tumor vector \vec{T} into the *Normal component* its linear models fit $Nc.\vec{T}$ onto the *Healthy State Model* and the *Disease component* $Dc.\vec{T}$ vector of residuals.

ease component—deviation from health vector—rather than the original data vector, several things are accomplished: (i) we emphasize the degree to which diseased tissue data are aberrant from healthy tissue data; (ii) we allow for a wide variability within the normal range; and (iii) we incorporate controls into the analysis. Working with the disease component of data has been shown both to outperform the use of original data and to bring out unique biology. Unlike direct comparison between normal and neoplastic tissue data, *DSGA* highlights the extent to which gene expression in a tumor is aberrant, whereas direct comparison tends to emphasize the background molecular signature of the progenitor cell type of the tumor. As we explain below, when combining the *DSGA* transformation with *Mapper*, we use as data the disease component of these data. We additionally define the guiding *Mapper* filter functions from the *DSGA* method.

1.3. Progression Analysis of Disease. We show now how to apply *Mapper* to *DSGA*-transformed data, with filter functions derived from the *DSGA* transformation. Importantly, the output of the procedure is a graph that highlights the core geometric shape of the data set of patients. As demonstrated in the next section, applying *PAD* to genomic data produces biologically meaningful insights and brings to light unique aspects of the biology of these tumors.

We begin with a data matrix from diseased tissue, in which columns are patients and rows are any genomic variable type, for example transcriptional microarray data. We assume we have tumor data vectors $\vec{T}_1, \vec{T}_2, \dots, \vec{T}_m$ and normal tissue data vectors $\vec{N}_1, \vec{N}_2, \dots, \vec{N}_k$ comprising the columns of the data matrix.

Step 1.

DSGA-transform all of the data and construct the following two matrices: (i) *Dc.mat*, the matrix whose columns $Dc.\vec{T}_1, \dots, Dc.\vec{T}_m$ are the disease components of the original tumor vectors $\vec{T}_1, \dots, \vec{T}_m$; (ii) *L1.mat*, a matrix whose columns $L1.\vec{N}_1, \dots, L1.\vec{N}_k$ are leave-one-out estimates of the deviation from healthy state by normal tissue data. Note that the columns of *L1.mat* constitute an estimate of the disease component of normal tissue. (iii) *L1Dc.mat*, the concatenated matrix with normal and tumor columns $L1.\vec{N}_1, \dots, L1.\vec{N}_k, Dc.\vec{T}_1, \dots, Dc.\vec{T}_m$.

Step 2.

Threshold data coordinates (genes, proteins, etc.) so that only the genes that show a significant deviation from the healthy state are retained in the data matrix from step 1. Any appropriate test for significance can be used.

Step 3.

Define *Mapper* filter functions on the data along which to perform the *Mapper* collapse to a graph. These functions should capture a biologically meaningful characteristic of the data. Essentially the data points are the individual columns of the *DSGA*-transformed data matrix, and for the filter functions we compute the vector magnitude in the L^p norm, as well as k powers of this magnitude. Below $f_{p,k}$ denotes the filter function, and \vec{V} denotes the column vector, either $Dc.\vec{T}_i$ or $L1.\vec{N}_j$. The coordinates are individual genes: $\vec{V} = \langle g_1, g_2, \dots, g_s \rangle$.

$$f_{p,k}(\vec{V}) = [\Sigma |g_r|^p]^{k/p}. \quad [2]$$

Note that if $k = 1$ and $p = 2$, the function simply computes the standard (Euclidean) vector magnitude of each column. Essentially, all these different filter functions, $f_{p,k}$, measure the overall amount of deviation from the null hypothesis, which is the *HSM*. Roughly, $f_{p,k}(Dc.\vec{T}_i)$ is large when a large number of genes deviates a lot from normal levels (the *HSM*) either in the positive direction (overexpression relative to normal) or the negative direction (underexpression relative to normal). Therefore, by using a variety of distance measurements, all these functions measure the extent to which a diseased tissue is different from

normal tissue. A tissue sample that has many genes exhibiting either increased or decreased activity relative to normal would show a large value of the filter $f_{p,k}$. A sample that resembles normal tissue in its gene activity will show a small value of $f_{p,k}$, close to 0. The effect of the different choices of p determining the choice of L^p norm is that, for larger values of p the weight of genes with larger expression levels is greater. Thus, the choice of p acts as an additional smooth threshold of genes.

Step 4.

Apply *Mapper* to the data obtained in step 2, using the filter functions defined in step 3. *Mapper* also requires that we define a distance function on the data: a measure of similarity between individual data points. The distance function used is the correlation distance.

2. Application of *PAD* to Breast Cancer Microarray Data

We applied the steps defined in the previous section to a breast cancer microarray gene expression data set (6). Normal tissue data were a set of 13 microarrays (4): four from reduction mammoplasty and nine normal tissue samples from cancer patients. Details of this analysis can be found in *SI Text*. The *DSGA* transformation and gene thresholding (steps 1 and 2) produced a data matrix with 262 rows (genes). *Mapper* filter functions were computed for the following parameters: k powers of the L^p distance with $p = 1, \dots, 5$ and powers $k = 1, \dots, 10$. Fig. 3 shows the output of *PAD* analysis for $p = 2$ and $k = 4$. Each node is a bin of tumors, and its color encodes the value of the filter function averaged across all of the data points in the bin, with blue denoting a low value and red encoding a large value. Thus, bins that are blue contain tumors whose expression is close to normal, whereas bins that are red contain tumors that generally have large deviation from normal along multiple genes, in both the positive and the negative direction. There are several groups of tumors that stand out. Basal tumors occupy most of the bins in the tumor sequence denoted as *ER*[−] sequence. They are immediately visible and stand out with large value (red) in the filter function: overall deviation from normal. Normal tissue samples all fall in the same bin together with 15 additional *ER*⁺ tumors. These are colored blue and show minimal overall deviation from normal according to the filter function. The known group of *her2*⁺ tumors is not yet visible, owing to the well-understood problem that only a small number of genes (on 17q) identify it, making them mathematically less visible, despite the fact that the small number of coordinates (17q genes) are biologically important. This discrepancy between mathematical and biological significance will be addressed in a later article. An additional long tumor sequence on the graph, the *ER*⁺ sequence showing large deviation from normal, is visible, as defined by the filter. This tumor sequence also consists of *ER*⁺ tumors, but unlike the first (blue) group of tumors, these are distinct from normal tissue in that the value of the filter function—the L^p magnitudes of the tumor vectors $Dc.\vec{T}_i$ in these bins—is very large. The breakdown of genes that most deviate from normal within the *ER*⁺ sequence tumors is given below in sections 2.4 and 2.5, but much of the positive gene activity centers on *Estrogen Receptor* and *c-MYB*. A subgroup of tumor bins is flanked by areas of sparse bins and is termed *c-MYB*⁺ tumors, because, as we show later in section 2.5, the list of significant genes points to crucial involvement of this and related genes. The *c-MYB*⁺ subset of tumors was also chosen to be the most dense segment of the *ER*⁺ sequence because it remains in the *PAD* output even when small bins containing only one data point are thresholded from the graph. This is very helpful to consider, because dropping the smallest bins provides a schematic of the denser part of data and corresponds to removing outliers. The simplified *PAD* output with small bins removed can be seen in *SI Text*. For the remainder of this section we analyze properties of these two very different subsets of *ER*⁺ tumors.

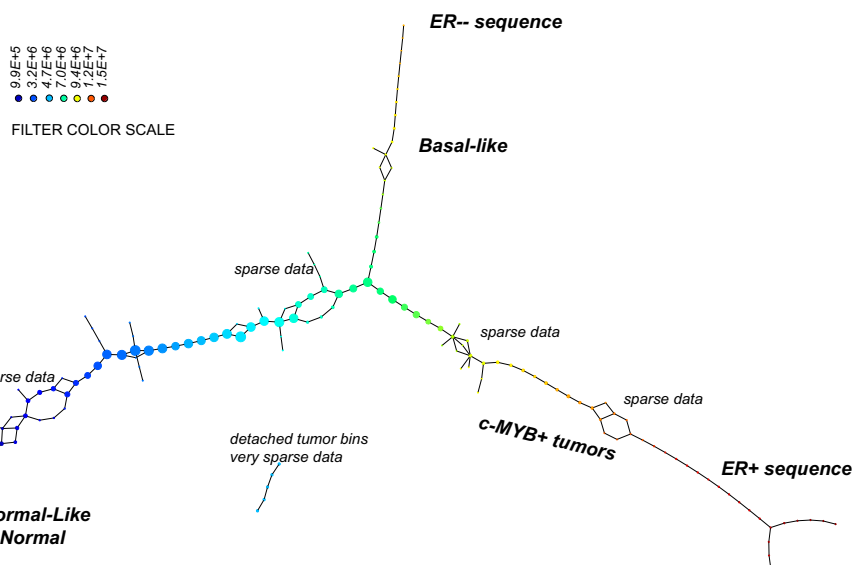


Fig. 3. PAD analysis of the *NK1* data. The output has three progression arms, because tumors (data points) are ordered by the magnitude of deviation from normal (the *HSM*). Each bin is colored by the mean of the filter map on the points. Blue bins contain tumors whose total deviation from *HSM* is small (normal and *Normal-like* tumors). Red bins contain tumors whose deviation from *HSM* is large. The image of *f* was subdivided into 15 intervals with 80% overlap. All bins are seen (outliers included). Regions of sparse data show branching. Several bins are disconnected from the main graph. The *ER*[−] arm consists mostly of *Basal* tumors. The *c-MYB*⁺ group was chosen within the *ER* arm as the tightest subset, between the two sparse regions.

The *Normal-like* (blue) group of tumors (15 tumors) constitutes 5% of the cohort. The low value of the filter function indicates little activity different from normal.

The $c\text{-MYB}^+$ (red) group of tumors (22 tumors) constitutes 7.5% of the cohort, or the more compact subset (outliers removed 14 tumors) 5% of ER^+ tumors. The high value of the filter function identifies these tumors as among the most distinct from normal tissue, showing extremely high activity in some gene groups (ER^+ , $c\text{-MYB}^+$) and low activity in others (innate immune genes), relative to normal tissue. This extreme deviation from normal molecular profiles, together with the biology of the overly active gene groups, and the excellent overall survival suggests that these tumors have a mechanism to respond in a protective way, antagonizing the presence of neoplastic tissue. In the next paragraphs we give evidence for the following two points: (i) $c\text{-MYB}^+$ breast cancer warrants being identified as a breast cancer group because it shows uniformity in molecular signature and clinical and survival properties, and because it is validated in other cancer data sets; and (ii) $c\text{-MYB}^+$ breast cancer is a unique group that does not fit into previously identified breast cancer types.

2.1. Survival Analysis. Survival analysis was performed on each of the two groups of ER^+ tumors: the blue *Normal-like* group and the red group that shows altered transcriptional activity in a large number of genes compared with the normal tissue, $c\text{-MYB}^+$ red group. Each group showed 100% overall survival, with no recurrence and no death from disease. Median time to follow-up was 10 y for the *Normal-like* group and 8.5 y for the $c\text{-MYB}^+$ tumors. It is important to note that survival information was not incorporated in the *DSGA* decomposition or the *Mapper* progression. We simply tested survival of groups of tumors that our *PAD* analysis found to stand out, purely on the basis of our two-step analysis: (i) *DSGA*, highlighting the distinction between normal and disease data, and (ii) *Mapper*, identifying subtle aspects in the shape of the data.

2.2. Comparison with Cluster Analysis Applied to the Same Data Matrix.

The *Normal-like* tumor group (blue) is often observed

through this type of analysis. However, the other group, *c-MYB*⁺ tumor group, was scattered across several clusters, as seen in Fig. 4. Thus, unlike *PAD*, cluster analysis was unable to identify this new group of tumors. This shows that the appearance of the new group of tumors was not due to the way data were transformed via *DSGA* nor to the specific method used for thresholding genes, but rather to the ability of *PAD* to identify subtle shape characteristics of the data set. Cluster analysis scattered the tumors in the *ER*⁺ tumor progression and even the very tight *c-MYB*⁺ tumor group. That the tumors in this group (22 in all, 14 without outliers) ought indeed to appear together is seen below, in sections 2.4–2.6, which show that the molecular signatures of these tumors are indeed very similar to one another and significantly distinct from other tumors.

2.3. Comparison with Molecular Subtype Classification. The 22 tumors in the *c-MYB*⁺ group were analyzed for molecular subtype (*Basal*, *ERBB2*, *Luminal A*, *Luminal B*, and *Normal-like*) (7) as previously assigned (6). Of the 22 tumors, only six had correlation >0.1 to one of the five centroids, the rest having been left unclassified. Five were classified as *Luminal A* and one as *Normal-like*. The rest of the *c-MYB*⁺ tumors were partially classified by the centroid they were closest to as follows: seven *Normal-Like*, six *Luminal A*, and three *Luminal B*. These assignments to subtype have changed (9) to be two *Normal-Like*, two *Luminal B*, and 18 *Luminal A*. This new assignment changes the subtype of 77% of tumors (17 of the 22 tumors have different assignment from their original one).

2.4. Prediction Analysis of Microarrays (PAM). *PAM* (10) was performed on *DSGA*-transformed data, using all genes, before thresholding (step 1 only). We wanted to investigate whether the two tumor groups, *c-MYB*⁺ and *Normal-like*, are good candidates for being molecular subtypes as far as their gene expression data were concerned. Using *PAM*, we wanted to determine whether they are (i) distinct from normal tissue, (ii) distinct from each other, and (iii) uniform within each group of tumors. Thus, we tested how successful *PAM* was in finding predictor variables for distinguishing these groups. The distinctions had extremely good

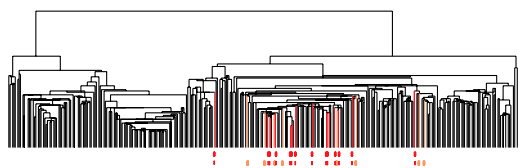


Fig. 4. Clustering vs. *PAD*. Can *Mapper* extract something new from the data that clustering does not? We compare the outputs of clustering (average linkage) vs. *Mapper* as applied to the same exact data matrix (*DSGA*-transformed *NKI*) to show that these two procedures are different. The bins defining the *c-MYB*⁺ group were marked on the cluster dendrogram (red for the tighter—no outliers—group, and orange for the larger *c-MYB*⁺ group containing outliers). The *c-MYB*⁺ tumors are scattered among different clusters, but *PAD* has been able to extract this group that turns out to be both statistically and biologically/clinically coherent.

error rates attained with very small numbers of genes, indicating that these groups of tumors satisfy all three conditions above. The output of the *PAM* analysis is found in *SI Text*. The distinction between *c-MYB*⁺ and normal is of particular interest: two predictor genes were able to distinguish between *c-MYB*⁺ group and normal tissue with error = 0. These predictor genes are *TSH-releasing hormone*, *TRH*, and *proprotein convertase subtilisin/kexin type 1*, *PCSK1*. Although it is important to remember that predictor variables need not be the most revealing about the underlying biology of the tumors, the fact that we are able to distinguish between *c-MYB*⁺ and normal with 0 error rate using only two genes is a strong indication that *c-MYB*⁺ is both significantly distinct from normal and significantly homogeneous as a class.

2.5. Significance of the Analysis of Microarrays (SAM). *SAM* (11) was performed on groups of tumors. Of special interest are the genes that are significantly different between (i) the *c-MYB*⁺ group and normal samples and (ii) the *c-MYB*⁺ group and the rest of the *ER*⁺ sequence in the *PAD* output. *Tables S1* and *S2* show the top genes in the output of these *SAM* analyses and demonstrate a significant set of differences between groups, as indicated by these lists of genes.

2.6. Testing the *c-MYB* Signature in the *c-MYB*⁺ Tumor Group. The *SAM* analysis identified the *c-MYB* gene to be among the significant top overexpressing genes (sixfold to 20-fold) in the *c-MYB*⁺ tumor group, both relative to normal tissue and relative to the rest of the *ER*⁺ tumor sequence in the *PAD* output. We wanted to find out whether other genes, known to be associated with (or downstream of) *c-MYB* overexpression (12), also show similar association in the *c-MYB*⁺ tumor group. We compared expression levels of known *c-MYB*-associated genes and computed *P* values using Student's *t* test; the results are found in *Table S3*. We tested the original rather than disease component values for the *c-MYB* signature. None of the genes listed as repressed by *MYB* overexpression showed significant reduction, but of the 45 genes listed as activated and present in the *Nederlands Kanker Instituut (NKI)* data, more than half (25 genes) had a *P* value <0.05 when values in the *c-MYB*⁺ group were compared with values in the normal group data.

2.7. Validation in Independent Breast Cancer Data. We validated the presence of the *c-MYB*⁺ group of tumors in two other breast cancer data sets: *Ullevål University Hospital (ULL)* (8) of 80 breast cancers, of which 52 were of ductal histological types, as were the *NKI* tumors and *HERSCH* (13) set of 232 tumors, of which 188 were primary breast tumors with good-quality RNA. We found the subset that best resembled the *c-MYB*⁺ among the identified *SAM* genes. Specifically, we considered *DSGA*-transformed tumor data along the 262 genes identified as *DSGA* sig-

nificant in the *NKI* data set, of which 255 genes were present in the *ULL* data and 221 in the *HERSCH* set. We further eliminated from the survival analysis step the tumors that had a very short follow-up time (<10 mo), as is standardly done because these short follow-up tumors affect negatively the reliability of survival analysis. Array mean-centered disease components were tested along the up and low sets of genes identified in the *SAM* analysis performed in the *NKI* data. Tumors were chosen on the basis of *SAM* genes in a two-step procedure: step 1 using two sets of *SAM* genes; step 2 using correlation along the 255 *DSGA* genes in common with the *ULL* set and the 221 *DSGA* genes in common with the *HERSCH* set. In step 1 we extracted tumors using two sets of *SAM* genes. First, we used the genes that were significant for the *PAD* progression arm *ER*⁺ sequence: the sequence of tumors leading up to the *c-MYB*⁺ group compared with normal, *Basal*, and *Normal-like* samples. Here we identified tumors which for at least 60% of the up *SAM* genes had expression levels higher than 33% of the tumors, and similarly, for 60% of the low *SAM* genes that had expression levels lower than 67% of the tumors. Second, we used the genes that were significantly distinct for the *c-MYB*⁺ subgroup compared with the rest of the tumors in *ER*⁺ sequence. This identified four tumors in the *ULL* set and 37 tumors in the *HERSCH* set. We then considered all of the tumors that were highly correlated (*r* > 0.68) to these top four tumors, along the 255 *DSGA* genes in the *ULL* set. Similarly, in the *HERSCH* set we identified tumors highly correlated (*r* > 0.60) to the top 37 tumors. This identified six tumors (13%) of the 46 total in *ULL* and 19 tumors (10%) of the total 188 in *HERSCH*. Finally, we tested survival in this group and again found them to have perfect survival and recurrence. Although this *c-MYB*⁺ subgroup consisted of only a few tumors, these constitute 13% of patients in *ULL* and 10% in *HERSCH*, thus higher than the 7.5% found in the first or *NKI* data set.

3. Discussion

We have introduced *PAD*, a method of analysis that takes into account the topology of data obtained from microarrays of disease tissue. First, *DSGA* highlights the expression pattern that deviates from normal (4). The second component of *PAD* consists in identifying the shape of *DSGA*-transformed data to access its topological properties beyond its cluster decomposition. Whereas cluster analysis identifies regions of higher density in these data, *Mapper* is able to find long gradual progressions, as is clearly demonstrated in this article. Here *PAD* identifies both quasi-parallel splits in progression, when a long string of data points suddenly splits into two gradually divergent progressions, as well as complete breaks, where data truly separate into disconnected regions. Moreover, *Mapper* creates a graph. This provides a means to visualize the shape of these data by way of a graph, and *Mapper* is flexible in the choice of guiding filter functions along which these data are collapsed to produce the graph. The filter functions are essentially a supervised step in the analysis, and different filter functions defined on the same data set highlight distinct shape features of these data. We note that *Mapper* is a much more general method to transform data into graphs, whereby filter functions can be chosen in a myriad possible ways. Different filter functions will highlight different aspects of the data. Indeed, several filter functions can be applied at once, thereby highlighting several aspects of the data at once. Moreover, owing in part to the simplicity of the graph output, the central problem of robustness of output can be addressed in a rigorous manner, using the concept of persistence (1). Thus, *Mapper*, in its complete generality, opens the door to study a wide range of data analysis problems. These and other aspects of *Mapper* will be discussed in further articles. Here we have attacked a very concrete type of omic data analysis problem, having defined the *Mapper* filter directly from the *DSGA* analysis as a measure of how aberrant the gene expression profile

Supporting Information

Nicolau et al. 10.1073/pnas.1102826108

SI Text

1. Microarray Data Analysis. We provide details for the microarray data analysis of the *Nederlands Kanker Instituut (NKI)* data (1) consisting of 295 tumors, the *Breast Cancer Normal (BCN)* data (2) consisting of 13 normal breast tissue samples, and the validation data sets *Ullevål University Hospital (ULL)* (3) consisting of 46 tumors of ductal histological type that had been in the study for longer than 10 mo and *HERSCH* (4) consisting of 188 primary breast tumors.

1.1. Data preprocessing. Data were retrieved, missing values imputed, then data were collapsed by UniGene cluster ID build 219, and genes present in both the tumor cohort and the normal data set were retained.

For *NKI*, data consisted of 24,479 GeneBank accession IDs on 295 tumor samples, all of which had at least 70% data. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were also transformed from the original \log_{10} values to \log_2 . Data were then collapsed (mean) by UniGene to the mean. The resulting data set consisted of 18,970 UniGene clusters.

For *BCN*, data from 13 normal tissue samples (nine nonneoplastic tissue from cancer patients, four reduction mammoplasty tissue) were retrieved with quality filters for each spot: (i) spot regression correlation $r > 0.6$, or (ii) channel 1 mean intensity/median background intensity > 1.5 , or (iii) channel 2 normalized (mean intensity/median background intensity) > 1.5 . Clones with 70% data were retained: 32,644 clone IDs. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were then collapsed by UniGene to 18,971 UniGene clusters. Of these, 12,237 UniGene IDs were in common with the *NKI* data set, and 17,441 were in common with the *ULL* data set (see below).

For *ULL*, data from 46 tumors were retrieved with quality filters for each spot: (i) spot regression correlation $r > 0.6$, or (ii) channel 1 mean intensity/median background intensity > 1.5 , or (i) channel 2 normalized (mean intensity/median background intensity) > 1.5 . Only clones with 70% good data were retained: 31,667 clone IDs. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were then combined with normal tissue data *BCN* and collapsed by UniGene to 17,441 UniGene clusters.

For *HERSCH*, data from 188 primary tumors were retrieved with quality filters for each spot: (i) spot regression correlation $r > 0.6$, or (ii) channel 1 mean intensity/median background intensity > 1.5 , or (iii) channel 2 normalized (mean intensity/median background intensity) > 1.5 . Only clones with 70% good data were retained: 32,644 clone IDs. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were then combined with normal tissue data *BCN* and collapsed by UniGene to 18,896 UniGene clusters.

1.2. Disease-Specific Genomic Analysis (DSGA). For *NKI* and *BCN*, data from tumors and normal tissue were combined along the common 12,237 UniGenes, and columns were normalized to have the magnitude of the mean vector magnitude of 13 normal tissue samples. The *Healthy State Model (HSM)* was constructed from normal tissue data $\{\vec{N}_1, \dots, \vec{N}_{13}\}$ as follows: *FLAT* construction (2) is a method to de-sparse the data in high dimensions by substituting for each normal tissue vector \vec{N}_i , its fit \hat{N}_i to a linear model in the other normal tissue vectors:

$$\hat{N}_i = \sum_{\substack{1 \leq j \leq 13 \\ j \neq i}} \beta_j \vec{N}_j.$$

This was shown to decrease noise in simulated data and help identify a good dimension reduction for *Principal Component*

Analysis (PCA). We use a method described in ref. 2 to compute the Wold invariant (6) designed to measure a version of signal-to-noise ratio:

$$W(l) = \left(\frac{\lambda_l^2}{\lambda_{l+1}^2 + \dots + \lambda_{13}^2} \right) \frac{(n-l-1)(13-l)}{(n+13-2l)}.$$

Fig. S1 plots $W(l)$ vs. the dimension l and shows a jump at $l = 10$, indicating that signal-to-noise ratio is higher at dimension 10, thereby justifying *PCA* dimension reduction of the *FLAT* normal data to 10. This produced the 10 dimensional *HSM*. Linear models are then used to compute the fitted tumor data matrix to the *HSM* (normal component *Nc.mat*) and the residuals (disease component *Dc.mat*). Along with tumor data, a leave-one-out procedure gives an estimate of the deviation of normal tissue data from the model of the healthy state *HSM*. Details of this procedure are found in ref. 2.

The validation data sets *ULL* and *HERSCH* were similarly transformed using the same normal data set *BCN*.

For gene thresholding, the 12,237 genes in the disease component matrix *Dc.mat* of tumors were reduced to 262 through the following method of testing for significance in deviation from the null hypothesis space. For each gene we computed the 5th and 95th percentiles of values in the disease components of the 295 tumors, and we recorded the larger of the two in absolute value and denoted the collection of these gene-by-gene deviations from normal by *MaxAbs595*. A histogram of these values is seen in Fig. S2. We then computed the 85th and 98th percentiles of *MaxAbs595* and denoted these as *relaxed* threshold and *stringent* threshold, respectively. A total of 1,836 genes exceeded the relaxed threshold, and 245 genes exceeded the stringent threshold. Genes were retained for further analysis if they passed the relaxed threshold and if they were also highly correlated ($r > 0.6$) to at least three genes that passed the stringent threshold. A total of 262 genes satisfied the condition. This method ensures that genes are retained in the analysis if they not only (i) deviate significantly from the null hypothesis space *HSM* but (ii) do so in groups of highly correlated genes. We denote the reduced matrix of disease component of *NKI* data: *nkiDc.mat*. The result of clustering the *nkiDc.mat* array and gene mean-centered can be found in supplementary folder Dataset S1: *nkiDc.AGmc.cdt*. It can be explored with *TreeView* (7), and all of the known clusters of genes can be observed, but because this is not germane to our present study we forgo any in-depth analysis of this clustering.

We did not follow the same thresholding procedure for the validation data sets *ULL* and *HERSCH*; rather, we found that of the 262 genes retained in the *NKI* data set, 255 genes were present in the *ULL* data and 221 in the *HERSCH* data.

1.3. Progression Analysis of Disease (PAD) on NKI. We give details of *PAD* on the reduced and *DSGA*-transformed *NKI* data matrix: *nkiDc.mat* of 295 tumors and 262 genes. First, this was combined with the leave-one-out matrix that estimates normal tissue: *bcnL1.mat*. The *Mapper* filter function was computed on each column vector, as explained in the main text (Eq. 2). The image space was then fragmented into 15 intervals, with 80% overlap. Two outputs of mapper were obtained: the first, which included all of the bins, can be found in Fig. 3 (main text). The second provides the tighter streamlined subset of *Mapper* output, by excluding all bins with only one data point in them. The two outputs appear side by side in Fig. S3.

1.4. Comparison with clustering. Although *Mapper* incorporates clustering at the local level, the final output captures a wide

range of characteristics that are obfuscated by the standard methods of clustering the entire data. We provide in Fig. S4 an expanded version of the comparison presented in Fig. 4 (main text) between clustering and *PAD* analysis, complete with heat maps and progressions of bins. It is important to note that the comparison is performed after both *Mapper* and clustering were applied to exactly the same data matrix. Thus, whatever transformations one might perform on the data, for example *DSGA*, and however genes are thresholded to provide a reduced number of genes used in the analysis, the final step of clustering vs. *Mapper* generates very different outputs. Because both clustering and *Mapper* are methods that identify the shape of the data, this comparison highlights the fact that shape characteristics identified by *Mapper* can be lost by clustering. Note as well that clustering has scattered the *c-MYB*⁺ tumor group among several clusters. This is a common problem known to clustering: data points will be segregated into separate clusters, and sometimes data points that are fairly close to one another will be torn apart and scattered into separate clusters. This is precisely what has happened with the *c-MYB*⁺ group. Despite how similar the *c-MYB*⁺ tumors are to one another, clustering has not kept them together.

2. Genes of Interest Analysis. We isolated a subgroup of tumors, *c-MYB*⁺, through the use of *PAD*. We provide *Prediction Analysis of Microarrays (PAM)* (8) analysis outputs for comparison of this group to the normal tissue group. We provide *Significance of the Analysis of Microarrays (SAM)* (9) analysis for genes most significantly distinct between the *c-MYB*⁺ group and the normal tissue group, as well as genes most significantly distinct between the *c-MYB*⁺ group and the most adjacent tumors to the *c-MYB*⁺ group in the *PAD* output, namely the tumors in the *ER*⁺ arm that are not part of the *c-MYB*⁺ group.

2.1. PAM. *PAM* finds a small set of *predictor genes* for distinguishing between two groups of tumors. Fig. S5 shows *PAM* output for comparing the *c-MYB*⁺ group to the normal tissue group.

2.2. SAM. *SAM* finds a large number of *significant genes* that behave differently between two groups of tumors. Table S1 shows *SAM* output genes significantly distinct between the *c-MYB*⁺ group and the rest of the *ER*⁺ arm of *PAD* output. Table S2 shows *SAM* output genes significantly distinct between the *c-MYB*⁺ group and the normal tissue group.

2.3. *c-MYB* signature. Genes that are believed to be downstream from *c-MYB* (10) were tested in the *c-MYB*⁺ group vs. normal tissue using a one-sided Student *t* test. Results are listed in Table S3.

1. van de Vijver MJ, et al. (2002) A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med* 347:1999–2009.
2. Nicolau M, Tibshirani R, Børresen-Dale AL, Jeffrey SS (2007) Disease-specific genomic analysis: identifying the signature of pathologic biology. *Bioinformatics* 23:957–965.
3. Langerød A, et al. (2007) TP53 mutation status and gene expression profiles are powerful prognostic markers of breast cancer. *Breast Cancer Res* 9:R30.
4. Herschkowitz JL, He X, Fan C, Perou CM (2008) The functional loss of the retinoblastoma tumour suppressor is a common event in basal-like and luminal B breast carcinomas. *Breast Cancer Res* 10:R75.
5. Troyanskaya O, et al. (2001) Missing value estimation methods for DNA microarrays. *Bioinformatics* 17:520–525.
6. Wold S, et al. (1978) Cross-validated estimation of the number of component in factor and principal components models. *Technometrics* 20:397–405.
7. Saldanha AJ (2004) Java Treeview—extensible visualization of microarray data. *Bioinformatics* 20:3246–3248.
8. Tibshirani R, Hastie T, Narasimhan B, Chu G (2002) Diagnosis of multiple cancer types by shrunk centroids of gene expression. *Proc Natl Acad Sci USA* 99:6567–6572.
9. Tusher VG, Tibshirani R, Chu G (2001) Significance analysis of microarrays applied to the ionizing radiation response. *Proc Natl Acad Sci USA* 98:5116–5121.
10. Ramsay RG, Gonda TJ (2008) MYB function in normal and cancer cells. *Nat Rev Cancer* 8:523–534.

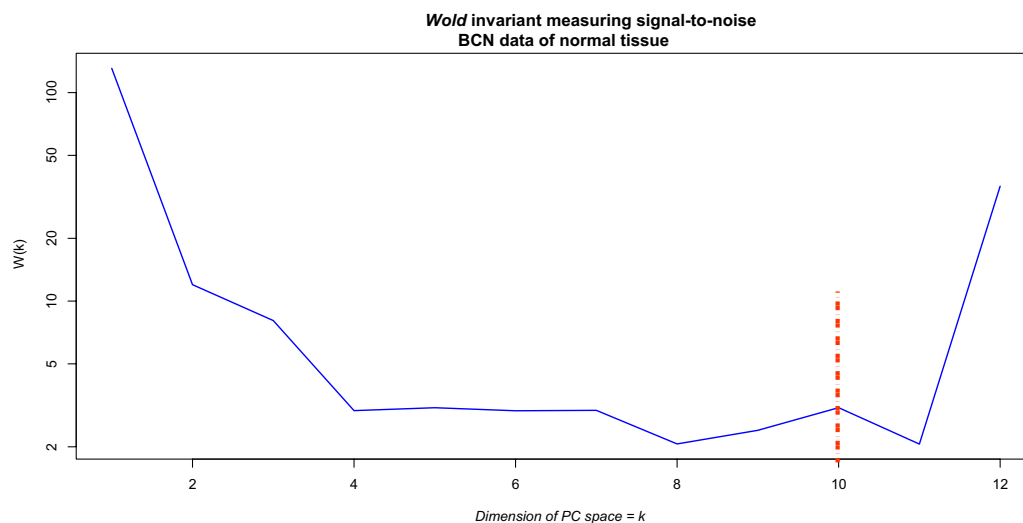


Fig. S1. The *wold* invariant is plotted as a function of the dimension reduction *K*. As the *wold* invariant is a measure of signal to noise, a local maximum in this plot indicates a good place to perform dimension reduction. In this case *K* = 10 is a good choice.

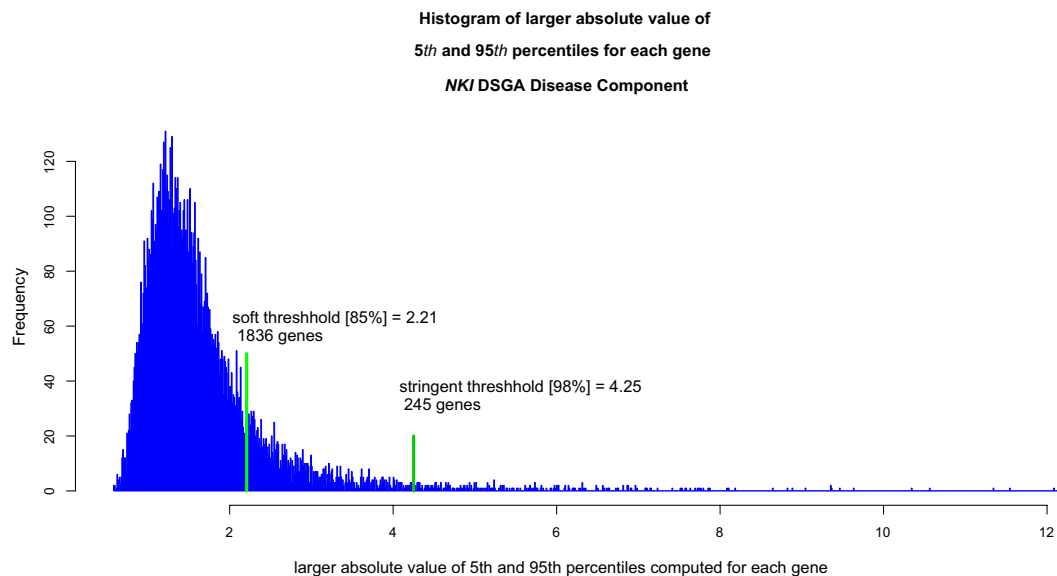


Fig. S2. For each gene, the 95th and 5th percentiles of expression levels in the disease component is computed. The larger of the two in absolute value denoted as Q_{gene} gives an estimate of the extent of deviation from normal for the gene. This deviation can be positive, indicating overexpression relative to normal levels, or negative, indicating underexpression relative to normal levels. The figure shows a histogram of the collection Q_{gene} of deviations from normal for the set of all genes. There are 1,836 genes for which this value exceeds the 85th percentile (*lax*-threshold genes) and 245 genes for which it exceeds the 95th percentile (*stringent*-threshold genes).

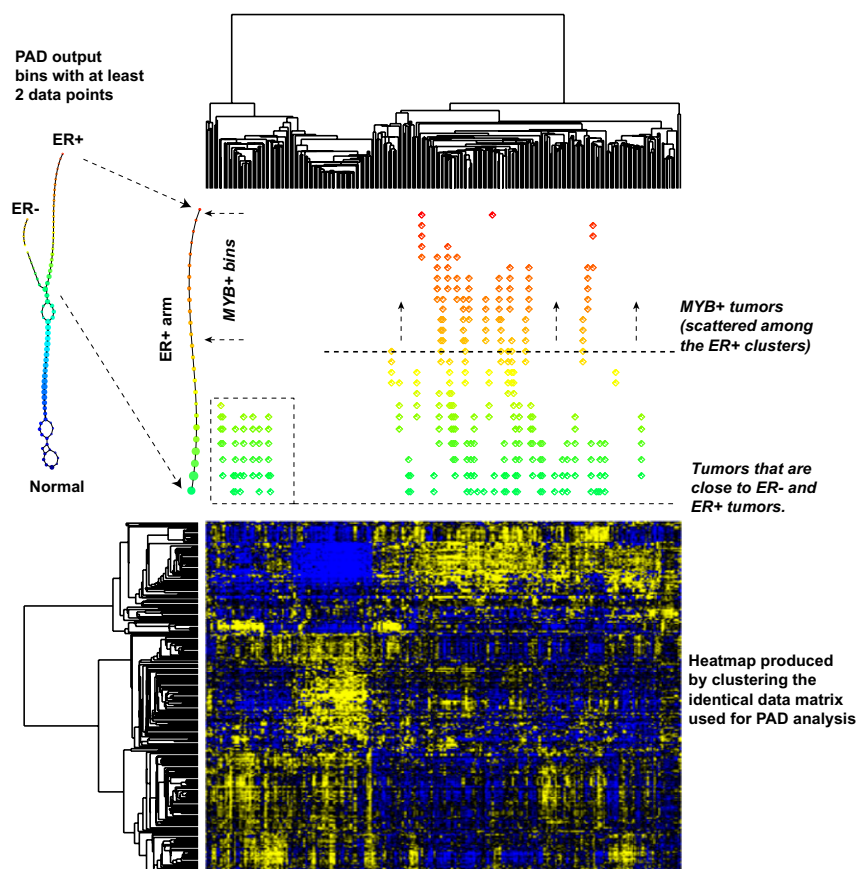


Fig. S4. Comparison between cluster analysis and PAD. Specifically, PAD consists of two major steps: the first step, *DSGA*, defines a transformation of the original data to detect extent of deviation from normal. It also provides a means to threshold genes so that only genes that deviate significantly from normal are retained. The second step, *Mapper*, involves detecting the shape of the data points in space. *Cluster analysis* is a different method to detect the shape of the data in space. This figure shows the difference between using cluster analysis as opposed to using *Mapper* to detect the shape of the same data matrix. We took the matrix whose columns are the disease components of the *DSGA*-transformed data, with only the 262 genes obtained by thresholding genes according to deviation from normal. This matrix was analyzed to detect its shape in space in two distinct ways: (i) it was clustered with associated heatmap and dendrograms shown, and (ii) it was processed with *Mapper*, with the output shown. The *ER*⁺ arm is magnified, and the position of each tumor in each consecutive bin is shown relative to its placement in the clustering dendrogram. It is easily visible that whereas the *c-MYB*⁺ group of tumors are close to one another in the PAD output, they are scattered throughout the *ER*⁺ portion of the clustering diagrams. It is important to note that the same matrix was fed into the *Mapper* and the cluster analysis. The figure shows these outputs to be very distinct. The figure does not and cannot identify which output is identifying features that deserve to be noticed: cluster analysis did not identify the *c-MYB*⁺ group, but it is not clear, simply on the basis of this figure, that the group is a real feature rather than an artifact of *Mapper*. It is through subsequent analysis methods that we see that the *c-MYB*⁺ group is indeed both mathematically and biologically distinct. Thus, the *PAM* analysis shows the group to be mathematically coherent and easily distinct, and functional exploration of the genes identified by *SAM* analysis, along with survival analysis of the group, show it to be a biologically coherent and meaningful group of tumors. This figure shows that the shape analysis provided by clustering is different from that provided by *Mapper*.

PAM analysis c-MYB+ group vs. Normal

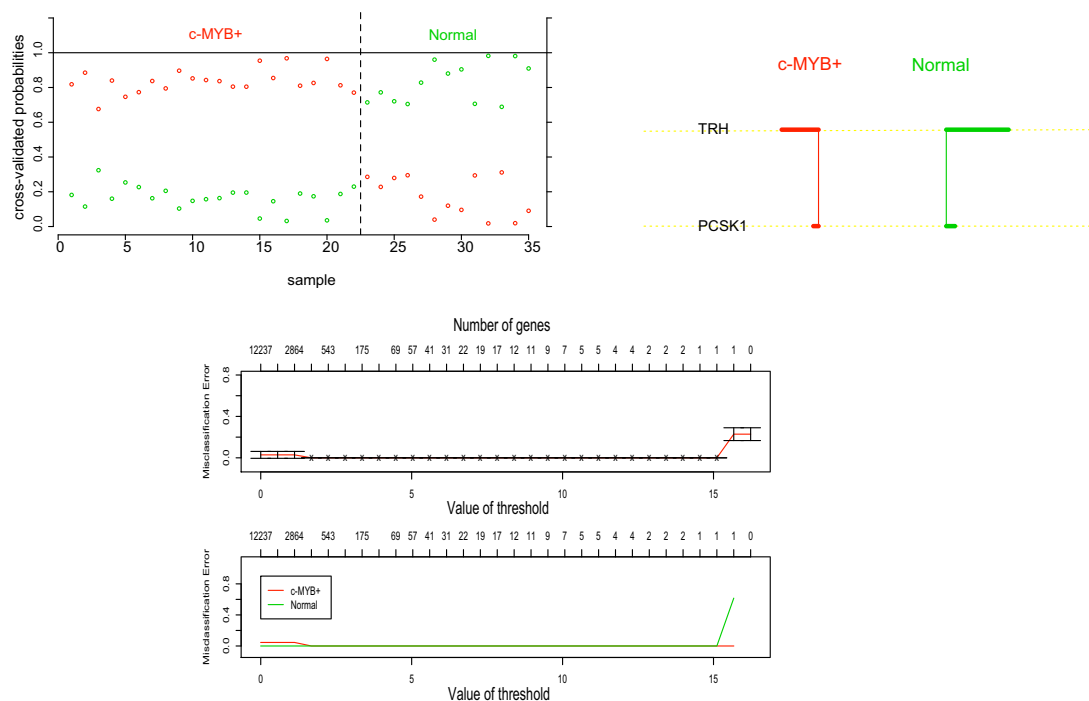


Fig. S5. Output of PAM analysis on the *c-MYB*⁺ group vs. *Normal* data. Two genes provide class prediction with error rate = 0: *TRH*, *TSH-releasing hormone*, and *PCSK1*, *proprotein convertase subtilisin kexin type 1*. The centroids, cross-validation probabilities, and misclassification error plots are shown.

Table S1. Genes significantly up-regulated and down-regulated in *MYB*⁺ vs. the rest of *ER*⁺ sequence

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. rest of ER ⁺ sequence
Hs.654446	MYB	0	MYB v-myb myeloblastosis viral oncogene homolog(avian)	Up
Hs.88417	SUSD3	0	SUSD3 Sushi domain containing 3	Up
Hs.414028	C9orf116	0	C9orf116 Chromosome 9 ORF 116	Up
Hs.532634	IFI27	5.15	IFI27 IFN, α -inducible protein 27 Hs.532634	Down
Hs.477891	CPB1	5.15	CPB1 Carboxypeptidase B1 (tissue) Hs.477891	Down
Hs.49760	ORC6L	5.15	ORC6L Origin recognition complex, subunit 6 like (yeast) Hs.49760	Down
Hs.517307	MX1	5.15	MX1 Myxovirus (influenza virus) resistance 1, IFN-inducible protein p78 (mouse) Hs.517307	Down
Hs.77367	CXCL9	5.15	CXCL9 Chemokine (C-X-C motif) ligand 9 Hs.77367	Down
Hs.501778	TRIM22	5.15	TRIM22 Tripartite motif-containing 22 Hs.501778	Down
Hs.521459	ADAMDEC1	5.15	ADAMDEC1 ADAM-like, decysin 1 Hs.521459	Down
Hs.458485	ISG15	5.15	ISG15 ISG15 ubiquitin-like modifier Hs.458485	Down
Hs.109225	VCAM1	5.15	VCAM1 Vascular cell adhesion molecule 1 Hs.109225	Down
Hs.17518	RSAD2	5.15	RSAD2 Radical S-adenosyl methionine domain containing 2 Hs.17518	Down
Hs.7155	CMPK2	5.15	CMPK2 Cytidine monophosphate (UMP-CMP) kinase 2, mitochondrial Hs.7155	Down
Hs.20315	IFIT1	6.51	IFIT1 IFN-induced protein with tetratricopeptide repeats 1 Hs.20315	Down
Hs.306777	GSDMB	6.51	GSDMB Gasdermin B Hs.306777	Down
Hs.715518	STAT1	6.51	STAT1 Signal transducer and activator of transcription 1, 91kDa Hs.715518	Down
Hs.709313	B2M	6.51	B2M Beta-2-microglobulin Hs.709313	Down
Hs.584823	PLA2G7	6.51	PLA2G7 Phospholipase A2, group VII (platelet-activating factor acetylhydrolase, plasma) Hs.584823	Down
Hs.181244	HLA-A	6.51	HLA-A Major histocompatibility complex, class I, A Hs.181244	Down
Hs.473341	SAMS1	6.51	SAMS1 SAM domain, SH3 domain and nuclear localization signals 1 Hs.473341	Down
Hs.523847	IFI6	6.51	IFI6 IFN, α -inducible protein 6 Hs.523847	Down
Hs.504641	CD163	6.51	CD163 CD163 molecule Hs.504641	Down
Hs.250615	CYP2A6	15.08	CYP2A6 Cytochrome P450, family 2, subfamily A, polypeptide 6 Hs.250615	Down
Hs.655652	LILRB2	15.08	LILRB2 Leukocyte Ig-like receptor, subfamily B (with TM and ITIM domains), member 2 Hs.655652	Down
Hs.459265	ISG20	15.08	ISG20 IFN stimulated exonuclease gene 20kDa Hs.459265	Down
Hs.926	MX2	15.08	MX2 Myxovirus (influenza virus) resistance 2 (mouse) Hs.926	Down
Hs.525157	TNFSF13B	15.08	TNFSF13B Tumor necrosis factor (ligand) superfamily, member 13b Hs.525157	Down
Hs.86859	GRB7	15.08	GRB7 Growth factor receptor-bound protein 7 Hs.86859	Down
Hs.352018	TAP1	15.08	TAP1 Transporter 1, ATP-binding cassette, subfamily B (MDR/TAP) Hs.352018	Down
Hs.32763	GRIA2	15.08	GRIA2 Glutamate receptor, ionotropic, AMPA 2 Hs.32763	Down
Hs.654585	PSMB9	15.08	PSMB9 Proteasome (prosome, macropain) subunit, β type, 9 (large multifunctional peptidase 2) Hs.654585	Down
Hs.718626	KIF20A	15.08	KIF20A Kinesin family member 20A Hs.718626	Down
Hs.474787	IL2RB	15.08	IL2RB Interleukin 2 receptor, β Hs.474787	Down
Hs.650174	HLA-E	15.08	HLA-E Major histocompatibility complex, class I, E Hs.650174	Down

Table S1. Cont.

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. rest of ER ⁺ sequence
Hs.143961	<i>CCL18</i>	15. 08	CCL18 Chemokine (C-C motif) ligand 18 (pulmonary and activation-regulated) Hs.143961	Down
Hs.81337	<i>LGALS9</i>	15. 08	LGALS9 Lectin, galactoside-binding, soluble, 9 Hs.81337	Down
Hs.474217	<i>CDC45L</i>	15. 08	CDC45L CDC45 cell division cycle 45-like (S. cerevisiae) Hs.474217	Down
Hs.301921	<i>CCR1</i>	15. 08	CCR1 Chemokine (C-C motif) receptor 1 Hs.301921	Down
Hs.16362	<i>P2RY6</i>	15. 08	P2RY6 Pyrimidinergic receptor P2Y, G protein coUpled, 6 Hs.16362	Down
Hs.419259	<i>REC8</i>	15. 08	REC8 REC8 homolog (yeast) Hs.419259	Down
Hs.591742	<i>IL7R</i>	15. 08	IL7R Interleukin 7 receptor Hs.591742	Down
Hs.647962	<i>ZIC1</i>	18.67	ZIC1 Zic family member 1 (odd-paired homolog, <i>Drosophila</i>) Hs.647962	Down
Hs.43388	<i>RTP4</i>	18. 67	RTP4 Receptor (chemosensory) transporter protein 4 Hs.43388	Down
Hs.376208	<i>LTB</i>	18. 67	LTB Lymphotoxin β (TNF sUperfamily, member 3) Hs.376208	Down
Hs.14623	<i>IFI30</i>	18. 67	IFI30 IFN, γ -inducible protein 30 Hs.14623	Down
Hs.660866	<i>CTSL2</i>	18. 67	CTSL2 Cathepsin L2 Hs.660866	Down
Hs.278658	<i>KRT86</i>	18. 67	KRT86 Keratin 86 Hs.278658	Down
Hs.1051	<i>GZMB</i>	18. 67	GZMB Granzyme B (granzyme 2, cytotoxic T lymphocyte-associated serine esterase 1) Hs.1051	Down
Hs.1594	<i>CENPA</i>	18. 67	CENPA Centromere protein A Hs.1594	Down
Hs.161985	<i>TMPRSS4</i>	18. 67	TMPRSS4 Transmembrane protease, serine 4 Hs.161985	Down
Hs.153752	<i>CDC25B</i>	18. 67	CDC25B Cell division cycle 25 homolog B (S. pombe) Hs.153752	Down
Hs.446352	<i>ERBB2</i>	18. 67	ERBB2 V-erb-b2 erythroblastic leukemia viral oncogene homolog 2, neuro/glioblastoma derived oncogene homolog (avian) Hs.446352	Down
Hs.497599	<i>WARS</i>	18. 67	WARS Tryptophanyl-tRNA synthetase Hs.497599	Down
Hs.182231	<i>TRH</i>	18. 67	TRH TSH-releasing hormone Hs.182231	Down
Hs.521903	<i>LY6E</i>	20.44	LY6E Lymphocyte antigen 6 complex, locus E Hs.521903	Down
Hs.370036	<i>CCR7</i>	20. 44	CCR7 Chemokine (C-C motif) receptor 7 Hs.370036	Down

Table S2. Genes significantly up-regulated and down-regulated in *MYB*+ vs. *Normal* tissue

UniGene build 219	Gene symbol	q value	Gene information	<i>MYB</i> level vs. normal
Hs.414028	<i>C9orf116</i>	0	C9orf116 Chromosome 9 ORF 116 Hs.414028	Up
Hs.406050	<i>DNALI1</i>	0	DNALI1 Dynein, axonemal, light intermediate chain 1 Hs.406050	Up
Hs.163484	<i>FOXA1</i>	0	FOXA1 Forkhead box A1 Hs.163484	Up
Hs.76704	[Hs.76704]	0	NA Transcribed locus Hs.76704	Up
Hs.654446	<i>MYB</i>	0	MYB V-myb myeloblastosis viral oncogene homolog (avian) Hs.654446	Up
Hs.88417	<i>SUSD3</i>	0	SUSD3 Sushi domain containing 3 Hs.88417	Up
Hs.494496	<i>FBP1</i>	0	FBP1 Fructose-1,6-bisphosphatase 1 Hs.494496	Up
Hs.448520	<i>SLC7A2</i>	0	SLC7A2 Solute carrier family 7 (cationic amino acid transporter, y+ system), member 2 Hs.448520	Up
Hs.534847	<i>C4A</i>	0	C4A Complement component 4A (Rodgers blood group) Hs.534847	Up
Hs.496240	<i>AR</i>	0	AR Androgen receptor Hs.496240	Up
Hs.631650	<i>GLT8D2</i>	0	GLT8D2 Glycosyltransferase 8 domain containing 2 Hs.631650	Up
Hs.91109	<i>PRR15</i>	0	PRR15 Proline rich 15 Hs.91109	Up
Hs.387057	<i>THSD4</i>	0	THSD4 Thrombospondin, type I, domain containing 4 Hs.387057	Up
Hs.98265	<i>ST6GAL2</i>	0	ST6GAL2 ST6 β -galactosamide α -2,6-sialyltransferase 2 Hs.98265	Up
Hs.208124	<i>ESR1</i>	0	ESR1 Estrogen receptor 1 Hs.208124	Up
Hs.111779	<i>SPARC</i>	0	SPARC Secreted protein, acidic, cysteine-rich (osteonectin) Hs.111779	Up
Hs.480819	<i>TBC1D9</i>	0	TBC1D9 TBC1 domain family, member 9 (with GRAM domain) Hs.480819	Up
Hs.437638	<i>XBP1</i>	0	XBP1 X-box binding protein 1 Hs.437638	Up
Hs.444414	<i>AFF3</i>	0	AFF3 AF4/FMR2 family, member 3 Hs.444414	Up
Hs.524134	<i>GATA3</i>	0	GATA3 GATA binding protein 3 Hs.524134	Up
Hs.467733	<i>GREB1</i>	0	GREB1 GREB1 protein Hs.467733	Up
Hs.458573	<i>PDGFRL</i>	0	PDGFRL Platelet-derived growth factor receptor-like Hs.458573	Up
Hs.210995	<i>CA12</i>	0	CA12 Carbonic anhydrase XII Hs.210995	Up
Hs.523468	<i>SCUBE2</i>	0	SCUBE2 Signal peptide, CUB domain, EGF-like 2 Hs.523468	Up
Hs.654370	<i>FAP</i>	0	FAP Fibroblast activation protein, α Hs.654370	Up
Hs.489142	<i>COL1A2</i>	0	COL1A2 Collagen, type I, α 2 Hs.489142	Up
Hs.416108	<i>CRKRS</i>	0	CRKRS Cdc2-related kinase, arginine/serine-rich Hs.416108	Up
Hs.371147	<i>THBS2</i>	0	THBS2 Thrombospondin 2 Hs.371147	Up
Hs.519601	<i>ID4</i>	0	ID4 Inhibitor of DNA binding 4, dominant negative helix-loop-helix protein Hs.519601	Up
Hs.100686	<i>AGR3</i>	0	AGR3 Anterior gradient homolog 3 (<i>Xenopus laevis</i>) Hs.100686	Up
Hs.435655	<i>ASPN</i>	0	ASPN Asporin Hs.435655	Up
Hs.425777	<i>UBE2L6</i>	0	UBE2L6 Ubiquitin-conjugating enzyme E2L 6 Hs.425777	Up
Hs.659093	[Hs.659093]	0	NA Transcribed locus Hs.659093	Up
Hs.93764	<i>CPA4</i>	0	CPA4 Carboxypeptidase A4 Hs.93764	Up
Hs.719277	<i>SLC39A6</i>	0	SLC39A6 Solute carrier family 39 (zinc transporter), member 6 Hs.719277	Up
Hs.604376	[Hs.604376]	0	NA Transcribed locus Hs.604376	Up
Hs.95612	<i>DSC2</i>	0	DSC2 Desmocollin 2 Hs.95612	Up
Hs.8059	<i>SYT4</i>	0	SYT4 Synaptotagmin IV Hs.8059	Up
Hs.1925	<i>DSG3</i>	0	DSG3 Desmoglein 3 (pemphigus vulgaris antigen) Hs.1925	Up
Hs.8786	<i>CHST2</i>	0	CHST2 Carbohydrate (<i>N</i> -acetylglucosamine-6-O) sulfotransferase 2 Hs.8786	Up
Hs.24950	<i>RGSS5</i>	0	RGSS5 Regulator of G protein signaling 5 Hs.24950	Up
Hs.19492	<i>PCDH8</i>	0	PCDH8 Protocadherin 8 Hs.19492	Up
Hs.520339	<i>COL10A1</i>	0	COL10A1 Collagen, type X, α 1 Hs.520339	Up
Hs.5210	<i>GMFG</i>	0.46	GMFG Glia maturation factor, γ Hs.5210	Up
Hs.497636	<i>LAMB3</i>	0.46	LAMB3 Laminin, β 3 Hs.497636	Up
Hs.6360	<i>TMCC2</i>	0.46	TMCC2 Transmembrane and coiled-coil domain family 2 Hs.6360	Up
Hs.34526	<i>CXCR6</i>	0.46	CXCR6 Chemokine (C-X-C motif) receptor 6 Hs.34526	Up
Hs.504115	<i>TRIM29</i>	0.85	TRIM29 Tripartite motif-containing 29 Hs.504115	Up

Table S2. Cont.

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. normal
Hs.1787	<i>PLP1</i>	0.85	PLP1 Proteolipid protein 1 Hs.1787	Up
Hs.523500	<i>CD2</i>	0.85	CD2 CD2 molecule Hs.523500	Up
Hs.131431	<i>EIF2AK2</i>	0.85	EIF2AK2 Eukaryotic translation initiation factor 2- α kinase 2 Hs.131431	Up
Hs.136348	<i>POSTN</i>	0.85	POSTN Periostin, osteoblast specific factor Hs.136348	Up
Hs.193235	<i>CPLX2</i>	0.85	CPLX2 Complexin 2 Hs.193235	Up
Hs.438	<i>MEOX1</i>	1.94	MEOX1 Mesenchyme homeobox 1 Hs.438	Up
Hs.405614	<i>CTHRC1</i>	1.94	CTHRC1 Collagen triple helix repeat containing 1 Hs.405614	Up
Hs.182231	<i>TRH</i>	0	TRH TSH-releasing hormone Hs.182231	Down
Hs.477891	<i>CPB1</i>	0	CPB1 Carboxypeptidase B1 (tissue) Hs.477891	Down
Hs.78977	<i>PCSK1</i>	0	PCSK1 Proprotein convertase subtilisin/kexin type 1 Hs.78977	Down
Hs.250615	<i>CYP2A6</i>	0	CYP2A6 Cytochrome P450, family 2, subfamily A, polypeptide 6 Hs.250615	Down
Hs.26770	<i>FABP7</i>	0	FABP7 Fatty acid binding protein 7, brain Hs.26770	Down
Hs.516874	<i>CHGB</i>	0	CHGB Chromogranin B (secretogranin 1) Hs.516874	Down
Hs.150793	<i>CHGA</i>	0	CHGA Chromogranin A (parathyroid secretory protein 1) Hs.150793	Down
Hs.77367	<i>CXCL9</i>	0	CXCL9 Chemokine (C-X-C motif) ligand 9 Hs.77367	Down
Hs.496843	<i>VGLL1</i>	0	VGLL1 Vestigial like 1 (<i>Drosophila</i>) Hs.496843	Down
Hs.268728	<i>TTYH1</i>	0	TTYH1 Tweety homolog 1 (<i>Drosophila</i>) Hs.268728	Down
Hs.416073	<i>S100A8</i>	0	S100A8 S100 calcium binding protein A8 Hs.416073	Down
Hs.473341	<i>SAMSN1</i>	0	SAMSN1 SAM domain, SH3 domain and nuclear localization signals 1 Hs.473341	Down
Hs.517307	<i>MX1</i>	0	MX1 Myxovirus (influenza virus) resistance 1, IFN-inducible protein p78 (mouse) Hs.517307	Down
Hs.532634	<i>IFI27</i>	0	IFI27 IFN, α -inducible protein 27 Hs.532634	Down
Hs.143961	<i>CCL18</i>	0	CCL18 Chemokine (C-C motif) ligand 18 (pulmonary and activation-regulated) Hs.143961	Down
Hs.458485	<i>ISG15</i>	0	ISG15 ISG15 ubiquitin-like modifier Hs.458485	Down
Hs.192859	<i>PCDH10</i>	0	PCDH10 Protocadherin 10 Hs.192859	Down
Hs.419259	<i>REC8</i>	0	REC8 REC8 homolog (yeast) Hs.419259	Down
Hs.470654	<i>CDCA7</i>	0	CDCA7 Cell division cycle associated 7 Hs.470654	Down
Hs.32763	<i>GRIA2</i>	0	GRIA2 Glutamate receptor, ionotropic, AMPA 2 Hs.32763	Down
Hs.415762	<i>LY6D</i>	0	LY6D Lymphocyte antigen 6 complex, locus D Hs.415762	Down
Hs.119689	<i>CGA</i>	0	CGA Glycoprotein hormones, α polypeptide Hs.119689	Down
Hs.278658	<i>KRT86</i>	0	KRT86 Keratin 86 Hs.278658	Down
Hs.17518	<i>RSAD2</i>	0	RSAD2 Radical S-adenosyl methionine domain containing 2 Hs.17518	Down
Hs.7155	<i>CMPK2</i>	0	CMPK2 Cytidine monophosphate (UMP-CMP) kinase 2, mitochondrial Hs.7155	Down
Hs.20315	<i>IFIT1</i>	0	IFIT1 IFN-induced protein with tetratricopeptide repeats 1 Hs.20315	Down
Hs.418167	<i>ALB</i>	0	ALB Albumin Hs.418167	Down
Hs.372578	<i>FAM65C</i>	0	FAM65C Family with sequence similarity 65, member C Hs.372578	Down
Hs.26225	<i>GABRP</i>	0	GABRP Gamma-aminobutyric acid (GABA) A receptor, ρ Hs.26225	Down
Hs.151254	<i>KLK7</i>	0	KLK7 Kallikrein-related peptidase 7 Hs.151254	Down
Hs.161985	<i>TMPRSS4</i>	0	TMPRSS4 Transmembrane protease, serine 4 Hs.161985	Down
Hs.376208	<i>LTB</i>	0	LTB Lymphotoxin β (TNF superfamily, member 3) Hs.376208	Down
Hs.414629	<i>CCL13</i>	0	CCL13 Chemokine (C-C motif) ligand 13 Hs.414629	Down
Hs.521459	<i>ADAMDEC1</i>	0	ADAMDEC1 ADAM-like, decysin 1 Hs.521459	Down
Hs.79361	<i>KLK6</i>	0	KLK6 Kallikrein-related peptidase 6 Hs.79361	Down
Hs.112405	<i>S100A9</i>	0	S100A9 S100 calcium binding protein A9 Hs.112405	Down
Hs.49760	<i>ORC6L</i>	0	ORC6L Origin recognition complex, subunit 6 like (yeast) Hs.49760	Down
Hs.647962	<i>ZIC1</i>	0	ZIC1 Zic family member 1 (odd-paired homolog, <i>Drosophila</i>) Hs.647962	Down
Hs.30743	<i>PRAME</i>	0	PRAME Preferentially expressed antigen in melanoma Hs.30743	Down

Table S2. Cont.

UniGene build 219	Gene symbol	<i>q</i> value	Gene information	<i>MYB</i> level vs. normal
Hs.2256	<i>MMP7</i>	0	MMP7 Matrix metalloproteinase 7 (matrilysin, uterine) Hs.2256	Down
Hs.523847	<i>IFI6</i>	0	IFI6 IFN, α -inducible protein 6 Hs.523847	Down
Hs.75285	<i>ITIH2</i>	0	ITIH2 Inter α (globulin) inhibitor H2 Hs.75285	Down
Hs.654550	<i>KRT13</i>	0	KRT13 Keratin 13 Hs.654550	Down
Hs.532635	<i>SERPINA6</i>	0	SERPINA6 Serpin peptidase inhibitor, clade A (α -1 antitrypsin), member 6 Hs.532635	Down
Hs.86859	<i>GRB7</i>	0.46	GRB7 Growth factor receptor-bound protein 7 Hs.86859	Down
Hs.514527	<i>BIRC5</i>	0.85	BIRC5 Baculoviral IAP repeat-containing 5 Hs.514527	Down
Hs.370036	<i>CCR7</i>	0.85	CCR7 Chemokine (C-C motif) receptor 7 Hs.370036	Down
Hs.22905	<i>RP13-102H20.1</i>	0.85	RP13-102H20.1 Hypothetical protein FLJ30058 Hs.22905	Down
Hs.660866	<i>CTSL2</i>	0.85	CTSL2 Cathepsin L2 Hs.660866	Down
Hs.315	<i>MUC2</i>	1.94	MUC2 Mucin 2, oligomeric mucus/gel-forming Hs.315	Down
Hs.63287	<i>CA9</i>	1.94	CA9 Carbonic anhydrase IX Hs.63287	Down
Hs.109225	<i>VCAM1</i>	1.94	VCAM1 Vascular cell adhesion molecule 1 Hs.109225	Down
Hs.501778	<i>TRIM22</i>	1.94	TRIM22 Tripartite motif-containing 22 Hs.501778	Down

Table S3. Testing the *MYB* signature genes

Gene symbol	Gene name	UniGene build 219	pval MYB*group_UP Normal_LO
<i>MYC</i>	V-myc myelocytomatosis viral oncogene homolog	Hs.202453	0.24
<i>MYB</i>	V-myb myeloblastosis viral oncogene homolog	Hs.654446	4.70E-05
<i>ADA</i>	Adenosine deaminase	Hs.654536	1.60E-10
<i>CDK1</i>	Cyclin-dependent kinase 1	Hs.334562	0.00019
<i>POLD1</i>	Polymerase (DNA directed), δ 1	Hs.279413	2.60E-11
<i>PRTN3</i>	Myeloblastin proteinase 3	Hs.928	0.00014
<i>CD4</i>	T-cell surface antigen T4/Leu-3	Hs.631659	1
<i>VEGF</i>	Vascular endothelial growth factor A	Hs.73793	0.62
<i>BCL2</i>	B-cell CLL/lymphoma 2	Hs.150749	0.97
<i>KIT</i>	Proto-oncogene c-Kit mast/stem cell growth factor receptor	Hs.479754	1
<i>CD34</i>	Hematopoietic progenitor cell antigen CD34	Hs.374990	1
<i>GATA3</i>	Transacting T-cell-specific transcription factor GATA-3	Hs.524134	0.00048
<i>MPO</i>	Myeloperoxidase	Hs.458272	0.012
<i>HSP70</i>	HSPA4 heat shock 70kDa protein 4	Hs.90093	0.00064
<i>H2A.Z</i>	H2AZ histone	Hs.119192	0.00028
<i>Adora2B</i>	Adenosine receptor 2B – chicken	Hs.167046	0.01
<i>Mcm4</i>	CDC21; CDC54; MGC33310; P1-CDC21; hCdc21	Hs.460184	2.20E-05
<i>GAS41</i>	YEATS4;Yeats domeain containing 4	Hs.4029	0.00078
<i>NMU</i>	Neuromedin U	Hs.418367	0.38
<i>CCNE1</i>	Cyclin E1	Hs.244723	0.00049
<i>CCNB1</i>	cyclin B1	Hs.23960	0.021
<i>CA1</i>	Carbonic anhydrase 1	Hs.23118	0.00037
<i>PDCD4</i>	Programmed cell death 4(neoplastic transformation inhibitor)	Hs.711490	0.019
<i>COL1A1</i>	Collagen type I, α 1	Hs.172928	1
<i>COL1A2</i>	Collagen type I, α 2	Hs.489142	1
<i>CD13 ANPEP</i>	Ananyl (membrane) animopeptidase	Hs.1239	0.96
<i>GBX2</i>	Gastrulation brain homeobox 2	Hs.184945	0.61
<i>Actn1</i>	Actinin, α 1	Hs.509765	0.9
<i>Birc3</i>	Baculoviral IAP repeat-containing 3	Hs.127799	1
<i>Casp6</i>	caspase 6, apoptosis-related cysteine peptidase	Hs.654616	3.60E-06
<i>Cbx4</i>	Chromobox homolog 4 (Pc class homolog, <i>Drosophila</i>)	Hs.714363	0.00073
<i>Copa</i>	coatomer protein complex, subunit α	Hs.162121	0.00017
<i>Hspa8</i>	Heat shock 70kDa protein 8	Hs.702021	1.80E-05
<i>Iqgap1</i>	IQ motif containing GTPase activating protein 1	Hs.430551	0.0047
<i>Lca CLTA</i>	Clathrin, light chain A	Hs.522114	9.00E-07
<i>Mad11l</i>	MAD1 mitotic arrest deficient-like 1 (yeast)	Hs.654838	7.30E-10
<i>Ppp3ca</i>	Protein phosphatase 3, catalytic subunit, α isozyme	Hs.435512	0.42
<i>SLC1A5</i>	Solute carrier family 1 (neutral amino acid transporter), member 5	Hs.631582	0.032
<i>Cox-2 PTGS2</i>	Prostaglandin-endoperoxide synthase 2 (prostaglandin G/H synthase and cyclooxygenase)	Hs.196384	0.28
<i>TCRd TRD@</i>	T-cell receptor δ locus	Hs.74647	0.79
<i>FABP5</i>	Fatty acid binding protein 5 (psoriasis-associated)	Hs.408061	1
<i>DHRS2</i>	Dehydrogenase/reductase (SDR family) member 2	Hs.272499	0.19
<i>TGFBI</i>	Transforming growth factor, β 1	Hs.645227	0.63
<i>CTNNA1</i>	Catenin (cadherin-associated protein), α -like 1	Hs.58488	0.00059