



Predicting Clinical Outcomes in Glioblastoma: An Application of Topological and Functional Data Analysis

Lorin Crawford, Anthea Monod, Andrew X. Chen, Sayan Mukherjee & Raúl Rabadán

To cite this article: Lorin Crawford, Anthea Monod, Andrew X. Chen, Sayan Mukherjee & Raúl Rabadán (2020) Predicting Clinical Outcomes in Glioblastoma: An Application of Topological and Functional Data Analysis, Journal of the American Statistical Association, 115:531, 1139-1150, DOI: [10.1080/01621459.2019.1671198](https://doi.org/10.1080/01621459.2019.1671198)

To link to this article: <https://doi.org/10.1080/01621459.2019.1671198>



View supplementary material [↗](#)



Published online: 17 Oct 2019.



Submit your article to this journal [↗](#)



Article views: 556



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 4 View citing articles [↗](#)



Predicting Clinical Outcomes in Glioblastoma: An Application of Topological and Functional Data Analysis

Lorin Crawford^{a,b,c}, Anthea Monod^d, Andrew X. Chen^e, Sayan Mukherjee^{f,g,h,i}, and Raúl Rabadán^e

^aDepartment of Biostatistics, Brown University, Providence, RI; ^bCenter for Statistical Sciences, Brown University, Providence, RI; ^cCenter for Computational Molecular Biology, Brown University, Providence, RI; ^dDepartment of Applied Mathematics, Tel Aviv University, Tel Aviv, Israel; ^eDepartment of Systems Biology, Columbia University, New York, NY; ^fDepartment of Statistical Science, Duke University, Durham, NC; ^gDepartment of Computer Science, Duke University, Durham, NC; ^hDepartment of Mathematics, Duke University, Durham, NC; ⁱDepartment of Bioinformatics & Biostatistics, Duke University, Durham, NC

ABSTRACT

Glioblastoma multiforme (GBM) is an aggressive form of human brain cancer that is under active study in the field of cancer biology. Its rapid progression and the relative time cost of obtaining molecular data make other readily available forms of data, such as images, an important resource for actionable measures in patients. Our goal is to use information given by medical images taken from GBM patients in statistical settings. To do this, we design a novel statistic—the smooth Euler characteristic transform (SECT)—that quantifies magnetic resonance images of tumors. Due to its well-defined inner product structure, the SECT can be used in a wider range of functional and nonparametric modeling approaches than other previously proposed topological summary statistics. When applied to a cohort of GBM patients, we find that the SECT is a better predictor of clinical outcomes than both existing tumor shape quantifications and common molecular assays. Specifically, we demonstrate that SECT features alone explain more of the variance in GBM patient survival than gene expression, volumetric features, and morphometric features. The main takeaways from our findings are thus 2-fold. First, they suggest that images contain valuable information that can play an important role in clinical prognosis and other medical decisions. Second, they show that the SECT is a viable tool for the broader study of medical imaging informatics. Supplementary materials for this article, including a standardized description of the materials available for reproducing the work, are available as an online supplement.

ARTICLE HISTORY

Received September 2017
Accepted September 2019

KEYWORDS

Euler characteristic;
Functional data;
Glioblastoma multiforme;
Shape statistics; Topological
data analysis

1. Introduction

The field of radiomics is focused on the extraction of quantitative features from medical magnetic resonance images (MRIs), typically constructed by tomography and digitally stored as shapes or surfaces. Quantifying geometric features from shapes in a way that is amenable to computational analyses has been a long-standing and fundamental challenge in both statistics and radiomics. Overcoming such a challenge would provide significant breakthroughs in broader scientific disciplines with the potential for real, practical impact. One particularly important application, where a viable quantification of shapes is needed, is the study of glioblastoma multiforme (GBM)—a glioma that materializes into aggressive, cancerous tumor growths within the human brain. GBM is a disease that is currently under active research in oncology; it is marked by characteristics that are not common in other cancers, such as spatial diffusivity and molecular heterogeneity. In human patients, it is a rapidly progressing disease with a post-diagnosis survival period of 12–15 months and, currently, there are only limited therapies available (Patel et al. 2014). Obtaining molecular information of GBM tumors entails an invasive medical procedure on the

patient that is costly in terms of both time and resources. In comparison, MRIs of these tumors are easily accessible and often readily available. Being able to effectively utilize MRIs of GBM tumors in computational settings increases the potential for well-developed statistical methodology to have a significant impact in cancer research and future treatment strategies.

There are two key aims of our work in this article: first, to quantify GBM tumor images to integrate medical imaging information into statistical models; and second, to explore the utility of medical imaging information in clinical studies of GBM. To achieve the first aim, we develop a novel statistic, the smooth Euler characteristic transform (SECT), that summarizes shape information of GBM MRIs as a collection of smooth curves. This allows the direct implementation of existing statistical models from functional data analysis (FDA); in particular, it allows tumor shape information to be used as a covariate in regression frameworks. To achieve the second aim, we study a cohort of individuals with publicly available MRIs from The Cancer Imaging Archive (TCIA) (Clark et al. 2013; Scarpace et al. 2016), as well as matched genomic and clinical data collected by The Cancer Genome Atlas (TCGA)

CONTACT Lorin Crawford  lorin_crawford@brown.edu  Department of Biostatistics, Brown University, Providence, RI; Anthea Monod  antheam@tauex.tau.ac.il

 Department of Applied Mathematics, School of Mathematical Sciences, Tel Aviv University, Tel Aviv 69978, Israel.

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/r/JASA.

 These materials were reviewed for reproducibility.

 Supplementary materials for this article are available online. Please go to www.tandfonline.com/r/JASA.

© 2019 American Statistical Association

(The Cancer Genome Atlas Research Network 2008). Through our extensive predictive analysis, we demonstrate a clinically relevant connection between the shape of brain malignancies and the variation of survival-based outcomes that are driven by molecular heterogeneity.

The remainder of this article is organized as follows. In Section 2, we outline the theoretical concepts used to quantify shape information of tumors and highlight their statistical utility; we also detail the construction of our statistic that summarizes tumor shape information, the SECT. In Section 3, we detail how regression methodologies for functional covariates are naturally suited to model the curves that capture tumor shape information. This connection with functional data allows us to specify a general regression model that intakes tumor shape information and turns out to be particularly powerful when conducting predictive inference. For our case study, we focus on Gaussian process (GP) regression with Markov chain Monte Carlo (MCMC) inference. In Section 4, we use the GP modeling framework to predict the clinical outcomes of GBM patients using gene expression data, existing morphometric and volumetric tumor image quantifications, and our proposed tumor shape summaries. Here, we perform a comparative study between each covariate type across different regressions generated by various covariance functions. Finally, in Section 5, we close with a discussion on possible future research.

2. Quantifying Tumor Images Using Topology

In this section, we develop a summary statistic that captures shape information from MRI images of GBM tumors, which will then be used as covariates in a regression model. The key strategy is to construct these statistics as a function that maps shapes into a Hilbert space. This function has two important properties: (i) it is injective and (ii) it admits a well-defined inner product structure. Notably, the inner product structure allows us to adapt ideas from FDA to specify general regression models that use shape summary statistics as predictor variables.

2.1. Background on Summary Statistics for Shape Data

Classical approaches represent shapes as a collection of landmark points (Kendall 1984; Bookstein 1997; Dryden and Mardia 1998). This data representation was implemented partly due to the limited image processing technology of the time. Current imaging technologies have since greatly improved and now allow three-dimensional shapes to be represented as meshes, which are collections of vertices, edges, and faces. Figure S1 depicts an example of a mesh representation for a brain tumor and ventricles. Recently, methods have been developed to generate automated geometric morphometrics for mesh representations (Boyer et al. 2011; Lipman and Daubechies 2011; Al-Aifari, Daubechies, and Lipman 2013; Boyer et al. 2015). However, despite these advancements, both user-specified and automated landmark-based methods are known to suffer from structural errors when comparing shapes that are highly dissimilar. Some examples of structural errors include: inaccurate pairwise correspondences between landmarks, alignment problems between dissimilar shapes, and global inconsistency of pairwise mappings. These structural errors tend to accumulate as the number

of landmarks imposed on each shape increases, and a high number of these points is often required to accurately capture shape information (especially when analyzing diverse shapes) (Gao et al. 2018). Such complications generally make landmark-based approaches less attractive.

Most recently, an approach known as the persistent homology transform (PHT) was developed to comprehensively address issues induced by landmark-based methods, and to maintain robust quantification performance for highly dissimilar and nonisomorphic shapes (Turner, Mukherjee, and Boyer 2014). While the PHT allows for the comparison of shapes without requiring landmarks, it does so by producing a collection of persistence diagrams—multiscale topological summaries used extensively in topological data analysis (TDA). This is restrictive because the geometry of the resulting summary statistics does not allow for an inner product structure that is amenable to (generalized) functional data models (Turner et al. 2014). We propose the SECT because it builds upon the theory of the PHT, in that it also produces a topological summary statistic, but it is constructed to be able to integrate shape information in regression-based methods. This proves to be particularly useful in our case study on predicting clinical outcomes in GBM.

2.2. Homology and Persistence

We begin by developing an intuition for *persistent homology* (Edelsbrunner, Letscher, and Zomorodian 2000; Zomorodian and Carlsson 2005), which is a foundational concept in TDA. Briefly, persistent homology can be viewed as the data-analytic counterpart to *homology*—a theoretical concept from algebraic topology, where the goal is to study the shape of abstract mathematical objects, such as sets and spaces, by counting occurrences of geometric patterns. In homology, the geometric patterns of interest are holes: homology groups provide a mathematical language for describing and keeping track of holes of an abstract mathematical object. The motivation behind classical algebraic topology is to then use these holes to distinguish between or suggest similarities among different abstract mathematical objects. For a more detailed review and theoretical discussion of these concepts, see the supplementary materials.

2.2.1. Homology

Homology is particularly relevant to our application and case study in GBM. Intuitively, not only does it describe contrasting physical tumor characteristics, but it also implicitly captures some information about the stage of disease progression. For example, *necrosis* is a form of cell injury which results in the premature death of cells. *Multifocality* is a radiological observation where individual tumor cells separate from the main mass and disperse elsewhere within the brain. From an imaging perspective, necrotic regions show up as dark regions (or holes) within a tumor, while multifocal tumors appear as segregated masses. Examples of both necrosis and multifocality captured by MRI images are shown in Figure S2. It has been suggested that the more necrosis or multifocality there is in a GBM tumor, the more aggressive the disease (Barker et al. 1996; Liu et al. 2015). Applying homology to radiomic studies not only identifies such

phenomena, but also tracks the number of times they occur and therefore provides a notion of disease severity.

Homology is indexed by integers: the 0th-degree homology captures the number of connected components in the shape, the 1st-degree homology captures the number of loops, and the 2nd-degree homology captures the number of voids. In the context of our GBM application, degree 0 homology corresponds to tumor masses and lesions. Degrees 1 and 2 homology correspond to necrosis, depending on whether we are analyzing two-dimensional image slices from an MRI or the three-dimensional tumor as a whole.

Despite its intuitive description, computing homology can be challenging. To this end, it is often convenient to represent the shape as a discrete union of simple building blocks, “glued” together in a combinatorial fashion. An important example of such a building block is the *simplex*: simplices are skeletal elements that take the form of vertices, edges, triangles (faces), tetrahedra, and other higher dimensional structures. A *simplicial complex* K is a collection of simplices and represents the discretization of a shape or tumor. Meshes that represent three-dimensional shapes are particular examples of finite simplicial complexes (again see Figure S1). There are two key interests in discretizing shapes into simplicial complexes. First, there exist efficient algorithms to compute homology for such discretizations; and second, discretization is essential for applying these abstract concepts to real data, where any given dataset will necessarily be finite.

In this article, we use the notation $H_k(K)$ to denote the k th homology group for the simplicial complex K . This corresponds to the collection of the k -dimensional elements of the simplicial complex. For example, $H_0(K)$ corresponds to the collection of vertices of the simplicial complex or, equivalently, to the collection of connected components of the shape (e.g., the masses and lesions of a tumor).

2.2.2. Persistent Homology

Persistent homology applies homology to data by continuously tracking the evolution of homology in the data at different scales (or resolutions). It can thus be seen as a way to extract and summarize geometric information. In persistent homology, the

index s of a *filtration* tracks the homological evolution. A filtration is a collection of simplicial complexes $\{K_s\}$ where the index s induces totally ordered sets $K_i \subseteq K_j$ for $i < j$. As s increases, the sequence of simplicial complexes $\{K_s\}$ also changes and grows. In this way, the index s of the filtration $\{K_s\}$ tracks the scale according to which the “shape” of the data changes and grows. The shape information at each scale s is encoded by the homology groups $H_k(K_s)$ of the simplicial complex K_s . More specifically, H_0 corresponds to the vertices, H_1 corresponds to edges, and H_2 corresponds to the faces of the simplicial complex or discretized shape. An example of a filtration is depicted in Figure 1. Here, the index s corresponds to the value of height function (which depends on some variable x and is discussed in detail further below) in the vertical direction ν . We see the evolution of vertices, edges, and a face appearing sequentially with height. Higher order structures are revealed as s increases.

Computing persistent homology produces a collection of intervals for each degree of homology, where each interval represents a k -dimensional topological feature (e.g., a connected component, loop, or void for a general, three-dimensional shape) that is “born” at the parameter value given by the left endpoint of the interval, and “dies” at the value at the right endpoint. The length of the interval corresponds to how long the topological feature “lives,” or persists. In this article, we consider these intervals to be represented by a *persistence diagram*. Persistence diagrams treat the start and endpoints of each interval as an ordered pair, and displays them as plotted points on a plane where the x -axis corresponds to birth time and the y -axis is the death time. Thus, one can consider a persistence diagram as a collection of points on and above the diagonal, with the set of points on the diagonal having infinite multiplicity (and included for regularity conditions; see the supplementary materials for further detail).

2.2.3. Persistent Homology Transform

The PHT captures shape information by collecting persistence diagrams of all degrees of homology, for all possible orientations of the shape. More formally, for a d -dimensional shape, the PHT results in d -many persistence diagrams arising from height function filtrations over infinitely many direction vectors on

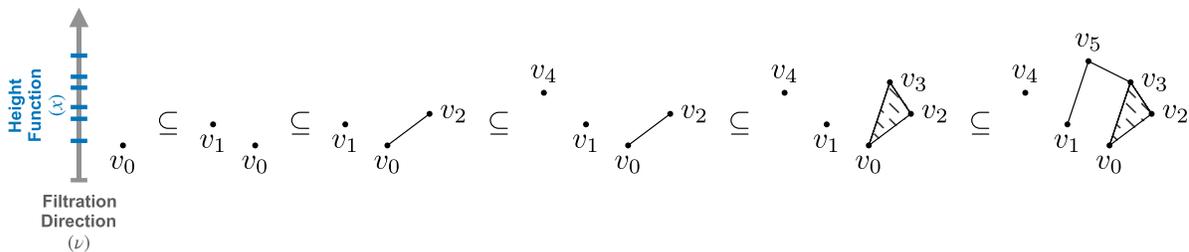


Figure 1. Demonstrating a filtration by height (as a function of x) in the vertical direction ν of a simplicial complex K . The inclusions (\subseteq ; from left to right) indicate the evolution of K with the filtration depicted by the hash marks (from the bottom to top) on the y -axis, where each element (i.e., vertex, edge, and face) is included at its maximal height x in the vertical direction ν , as given in Equation (1) from Section 2.3. Starting from the lefthand side, the vertex v_0 is first added and it remains present in each subsequent inclusion, while the vertex v_1 emerges during the second inclusion and then remains. Similarly, the edge v_0v_2 connecting vertices v_0 and v_2 appears at the third inclusion, and the face $\langle v_0v_1v_2 \rangle$ forms at the fourth. Altogether, according to these inclusions, the collection of vertices in K are shown to evolve as follows: (i) $\{v_0\}$; (ii) $\{v_0, v_1\}$; (iii) $\{v_0, v_1, v_2\}$; (iv) $\{v_0, v_1, v_2, v_4\}$; (v) $\{v_0, v_1, v_2, v_3, v_4\}$; and at the final inclusion, (vi) $\{v_0, v_1, v_2, v_3, v_4, v_5\}$. Similarly, the collection of edges in K evolves as follows: $\langle v_0v_2 \rangle$ appears at the third, $\langle v_0v_2, v_0v_3, v_2v_3 \rangle$ emerges at the fifth; and $\langle v_0v_2, v_0v_3, v_2v_3, v_1v_5, v_3v_5 \rangle$ forms at the final inclusion. This evolution of vertices, edges, and faces forms (persistence) vector spaces of i -chains, and may be written as such for concise notation (see the supplementary materials for further details). A version of this figure has been previously published (Turner, Mukherjee, and Boyer 2014).

the surface of the sphere. The space of persistence diagrams is a complicated, but theoretically well-defined probability space (Mileyko, Mukherjee, and Harer 2011). In particular, it is a metric space, meaning that distances between persistence diagrams may be defined. This is important because distances between PHT summary statistics provide a way of comparing shapes. The injectivity of the PHT for two- and three-dimensional shapes (Turner, Mukherjee, and Boyer 2014), or the one-to-one relation between the shape itself and its infinite collection of persistence diagrams, guarantees that the PHT effectively summarizes all relevant information about the shape.

Considering all possible directions on the surface of the sphere to summarize shape information is particularly well-suited to our radiomics application. MRI scans of the brain are known to be subject to noise: the positioning of patients' heads could vary both between patients and individual scans, causing image registration issues. Considering all directions on the surface of the sphere bypasses this problem, and incorporates perturbations directly into the statistic. This is an important feature of the PHT that we retain in the development of the SECT. We expand upon the PHT to produce a collection of continuous, piecewise linear functions that live in Hilbert space \mathbb{L}^2 . The corresponding inner product structure inherent to Hilbert spaces allows us to apply the SECT to a much broader set of statistical methodologies. It is worth noting that for select covariance functions, the PHT can be adapted to nonparametric statistical models (Kwitt et al. 2015; Reininghaus et al. 2015; Kusano, Fukumizu, and Hiraoka 2018), but this class is considerably limited.

2.3. Smooth Euler Characteristic Transform

While the SECT uses the same underlying mathematical principles as the PHT, it produces a collection of continuous, piecewise linear functions rather than persistence diagrams. The SECT implements persistent homology via the Euler characteristic (EC), which is a topological invariant that appears in many branches of mathematics. In terms of homology, the EC counts the ranks of the homology groups (i.e., the Betti numbers, β_k , for the k th homology group H_k) in an alternating sum and thus reduces the mathematical description of holes in a topological space from an algebraic group structure to an integer.

Definition 1. Let X be an arbitrary topological space, $H_k(X)$ be the k th homology group of X , and β_k be the rank of $H_k(X)$. The Euler characteristic (EC) $\chi(X)$ of X is the alternating sum

$$\chi(X) = \beta_0 - \beta_1 + \beta_2 - \beta_3 + \dots = \sum_{k=0}^{\infty} (-1)^k \beta_k.$$

For a discretized shape or surface in three dimensions represented as a simplicial complex K , the EC may be analogously defined by the number of simplices in K by

$$\chi(K) = V - E + F,$$

where V , E , and F are the numbers of vertices (0-simplices), edges (1-simplices), and faces (2-simplices), respectively.

Just as homology may be augmented to persistent homology by considering a filtration, ECs may also be calculated with

respect to a filtration. The result is an EC curve, which tracks the progression of the EC as a function with respect to the filtration. Let the dimension $d = \{2, 3\}$, and fix a direction ν on the surface of the unit circle or sphere S^{d-1} (where $\nu \in S^{d-1}$). Let \mathcal{M}_{d-1} be the set of all closed, compact subsets (shapes) embedded in \mathbb{R}^d that can be represented in a finite, discrete manner as simplicial complexes (Edelsbrunner and Harer 2010). Next, denote the simplicial complex representation of $M \in \mathcal{M}_{d-1}$ by K , and let K_ν indicate the ν -orientation of K . The *sublevel set filtration* of K_ν parameterized by a height function $r(\bullet, \bullet)$ is the set $\{x \in K : x \cdot \nu \leq r\}$. The ν -directional parameter height function $r_\nu(\bullet, \bullet)$ is

$$r : K \times S^{d-1} \rightarrow \mathbb{R} \\ \{x, \nu\} \mapsto x \cdot \nu. \quad (1)$$

Denote the extremal heights from this filtration by

$$a_\nu := \min\{r_\nu(x), x \in K\}, \\ b_\nu := \max\{r_\nu(x), x \in K\}.$$

We use the subscript notation to denote the simplicial complex representation K of a shape M , in the direction ν , as K_ν for $d = \{2, 3\}$. Similarly, we use the superscript notation K_ν^x to denote the varying simplicial complex of K_ν , generated by a sublevel set filtration with respect to Equation (1) and defined by varying $x \in K_\nu$.

Definition 2. The *EC curve* of K (which discretizes M) in the direction ν is defined by

$$\chi_\nu^K : [a_\nu, b_\nu] \rightarrow \mathbb{Z} \subset \mathbb{R} \\ x \mapsto \chi(K_\nu^x). \quad (2)$$

The EC curve tracks the evolution of the EC up to (and including) the largest subcomplex of K_ν^x contained in the sublevel set $r_\nu^{-1}((-\infty, x])$. See Figure 2(a) for an illustrative example of the evolution of the EC on the two-dimensional contour of a hand. Here, the direction is the horizontal direction to the right of the y -axis. The value of the EC changes as the sweep over the palm first reveals the thumb, and then the separation between the ring and pinky fingers, followed very shortly by the separation between the index and middle fingers, and so on. The EC curve of this filtration is plotted in Figure 2(b).

The same rotational summary technique of the PHT may be adapted to ECs as follows.

Definition 3 (Previously in Turner, Mukherjee, and Boyer (2014)).

In considering a directional sweep over the surface of the sphere S^{d-1} , and calculating the corresponding EC curves χ_ν^K of the finite simplicial complex representations K_ν for every direction $\nu \in S^{d-1}$, the *Euler characteristic transform (ECT)* is defined as follows:

$$\text{ECT}(K) : S^{d-1} \rightarrow \mathbb{Z}^{\mathbb{R}} \\ \nu \mapsto \chi(K_\nu). \quad (3)$$

In other words, the ECT of a shape collects EC curves of the shape, over all directions on the surface of the sphere.

The EC curve in Equation (2) and its corresponding ECT in Equation (3) are piecewise constant, integer-valued functions.

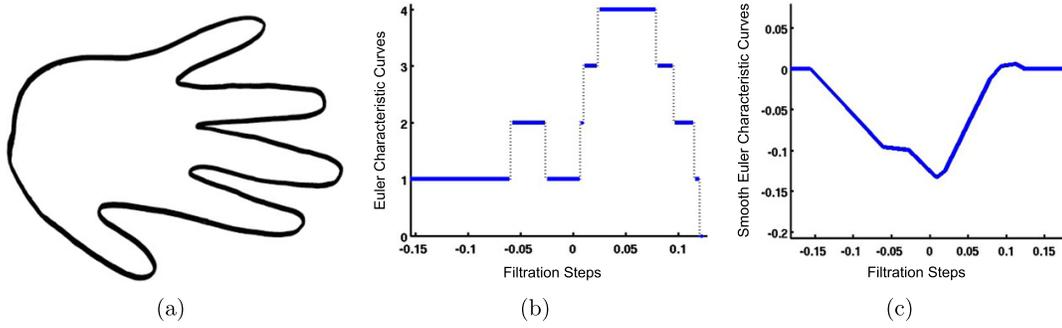


Figure 2. Example illustrating the evolutionary tracking of topological features for a given shape. (a) A two-dimensional contour of a hand, for which the Euler characteristic (EC) is calculated and tracked with respect to a horizontal filtration. (b) The Euler characteristic curve of the two-dimensional contour. At the leftmost filtration level on the x -axis, the EC value is equal to 1 on the y -axis. This indicates one connected component corresponding to the thumb in the hand. Just before level -0.05 , the index finger appears, which increases the value of the EC to 2. (c) The corresponding smooth Euler characteristic (SEC) curve. This found by taking the mean value of the EC curve in (b), subtracting it from every point along the x -axis, and integrating over the range of extremal heights (see Section 2.3). Since there are no holes in a hand, the EC reduces to the sum of connected components as they appear with the filtration. A version of this figure has been previously published (Turner, Mukherjee, and Boyer 2014).

These discontinuities can affect the stability of this representation (e.g., see Figure 2(b) where there are sharp jumps in the curve). Therefore, we propose a formulation of Equation (3) that allows for a type of summary that can be used in a wider range of statistical analyses. We do this by smoothing a centered variant of function. The centered variant is given by taking the mean of curve $\bar{\chi}_\nu^K$ over $[a_\nu, b_\nu]$ and subtracting it from the EC $\chi_\nu^K(x)$ at every $x \in [a_\nu, b_\nu]$. This produces a centered EC curve in the direction $\nu \in S^{d-1}$,

$$Z_\nu^K : [a_\nu, b_\nu] \rightarrow \mathbb{R} \quad (4)$$

$$x \mapsto \chi_\nu^K(x) - \bar{\chi}_\nu^K.$$

We set the value of Z_ν^K to be zero outside the interval $[a_\nu, b_\nu]$ by default. Integrating the curve gives the following smoothed construct.

Definition 4. The centered, cumulative Euler characteristic curve or smooth Euler characteristic curve (SEC), for a fixed direction $\nu \in S^{d-1}$, is defined for all $y \in \mathbb{R}$ as

$$SEC(K) : \mathbb{R} \rightarrow \mathbb{L}^2 \quad (5)$$

$$F_\nu^K(y) := \int_{-\infty}^y Z_\nu^K(x) dx.$$

The SEC is a continuous, piecewise linear function with compact support $[a_\nu, b_\nu]$ by construction. Therefore, it is an element of the Hilbert space \mathbb{L}^2 of square integrable functions on \mathbb{R} . The counterpart to Figure 2(b), smoothed by the procedure described above resulting in the SEC, is visually illustrated in Figure 2(c). We now formally define the SECT.

Definition 5. The SECT for a simplicial complex K of a shape $M \subset \mathbb{R}^d$, with $d = \{2, 3\}$, is the map

$$SECT(K) : S^{d-1} \rightarrow \mathbb{L}^2[a_\nu, b_\nu] \quad (6)$$

$$\nu \mapsto F_\nu^K(b_\nu)$$

for all $\nu \in S^{d-1}$. Each curve F_ν^K is also an element in the Hilbert space \mathbb{L}^2 . The following metric can therefore be used

to define distances between two simplicial complexes (discrete shape representations) K_1 and K_2 ,

$$\text{dist}_{\mathcal{M}_{d-1}}^{\text{SECT}}(K_1, K_2) := \left(\int_{S^{d-1}} \|F_\nu^{K_1} - F_\nu^{K_2}\|^2 d\nu \right)^{1/2}. \quad (7)$$

The advantage of the SECT over the PHT is that SECT summaries are a collection of curves and have a Hilbert space structure. This means that their structure allows for quantitative comparisons using the full scope of functional and nonparametric statistical methodology. The SECT is also an injective map and the following corollary is an immediate consequence of previous results (Turner, Mukherjee, and Boyer 2014).

Corollary 1. The SECT is injective for two- and three-dimensional shapes, that is, when the domain is \mathcal{M}_{d-1} for $d = \{2, 3\}$.

The injective property of the SECT suggests that it concisely summarizes the original shape data. Mathematically, the SECT maps between the space of all shapes with a finite simplicial complex representation \mathcal{M}_{d-1} and the Hilbert space \mathbb{L}^2 . Thus, injectivity between these two spaces means that for a given SECT statistic in \mathbb{L}^2 , there is a (unique) corresponding shape with some finite complex representation in \mathcal{M}_{d-1} . However, note that enough directions $\nu \in S^{d-1}$ must be taken for this corollary to hold since, for any one fixed direction, it is not true that the EC curve (upon which the SECT construction depends) is injective. An illustration of this fact is depicted in Figure S3. To determine the number of directions to use in practice, we perform a sensitivity analysis with many different combinations of numbers of directions and sublevel sets. In our application of interest and case study, we find prediction results (with SECT features as predictor variables) to be reasonably robust to our final choice of numerical parameters.

2.3.1. Small Note on Information Loss

Based on its homological definition and filtration (defined by a height function), the SECT will always capture all topological and integral geometric (i.e., size) information about a shape. However, there are instances where information about texture-

based features may not be captured in the SECT summary statistic. Intuitively, this is dependent upon the granularity of the filtration defined by the height function. In theory, a continuous height function would mitigate this issue, but this is difficult to implement in practice. As a result, coarse filtrations with too few sublevel sets will cause the SECT to miss or “step over” very local undulations in a shape. In the context of our GBM case study, this can occur if a particular tumor is made up of small focal lesions that are collectively important in explaining survival outcomes. For those cases, we would not observe this variance and presumably suffer in predictive performance.

3. Functional Regression Models With Tumor Shape Information as Covariates

In the previous section, we formally specified the SECT which allows us to map shapes into a space that (i) is represented by collection of curves and (ii) has a well-defined inner product structure. We will now discuss how FDA is a particularly suitable framework to specify a general regression model that uses tumor shape information (in the form of topological summary statistics, captured by the SECT) as covariates. The goal of FDA is to model data that are continuous functions (e.g., curves, response surfaces, or images) (Müller and Stadtmüller 2005; Ferraty and Vieu 2006; Ramsay 2006; Müller and Yao 2008; Pomann et al. 2016). The key idea here is that these functions can be considered as elements in a Hilbert space for which one can specify statistical models using stochastic processes (Preda 2007; Kadri et al. 2010; Morris 2015; Wang, Chiou, and Müller 2016). In this article, we will use a class of stochastic processes that is often referred to as GPs (Wahba 1997; Pillai et al. 2007; Yuan and Cai 2010).

3.1. Gaussian Process Regression

Denote the shape information (i.e., SECT representation) of a GBM MRI scan as $\mathbf{F}(t) = \{F_\nu\}_{\nu=1}^m$ measured over m directions. A functional linear model considers a continuous response variable \mathbf{y} and covariates that are square integrable functions $\mathbf{F}(t)$ on the real interval \mathcal{T} , where $t \in \mathcal{T}$ under the following parametric form (Müller and Stadtmüller 2005),

$$\mathbf{y} = \langle \mathbf{F}(t), \boldsymbol{\alpha}(t) \rangle + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}),$$

where the residual noise $\boldsymbol{\varepsilon}$ is assumed to follow a multivariate normal distribution with mean zero and scaled variance parameter τ^2 , and \mathbf{I} is used to denote the identity matrix. Notice that similar to traditional linear regression models, $\boldsymbol{\alpha}(t)$ is an unknown smooth parameter function that is now square integrable on the domain \mathcal{T} , and $\langle \bullet, \bullet \rangle$ denotes a well-defined inner product. In the context of our radiomic case study, the assumption of a linear relationship between the response variable \mathbf{y} and functional covariates $\mathbf{F}(t)$ may be too restrictive. For example, when modeling the topological landscape of brain tumors (as we will do in Section 4), it is reasonable to assume that interactions between modes of brain activity extend well beyond additivity (Friston et al. 2000). As a result, we formulate a general functional regression model that has the flexibility to incorporate possible nonlinear interactions. The methodology we use is GP regression.

There are two key characteristics of a GP regression model. The first key element is a positive definite covariance function, $\sigma : \mathbb{L}^2 \times \mathbb{L}^2 \rightarrow \mathbb{R}$, where again \mathbb{L}^2 is the Hilbert space of the SECT functional covariates such that $\mathbf{F}(t) \in \mathbb{L}^2$. The second key element is the reproducing kernel Hilbert space (RKHS) that is induced by the covariance function. Given the eigenfunctions $\{\psi_l\}_{l=1}^\infty$ and eigenvalues $\{\lambda_l\}_{l=1}^\infty$ of the finite integral operator defined by the covariance function (Mercer 1909), we have

$$\int_{\mathcal{T}} \sigma(\mathbf{u}, \mathbf{v}) d(\mathbf{u}, \mathbf{v}) < \infty, \quad \lambda_l \psi_l(\mathbf{u}) = \int_{\mathcal{T}} \sigma(\mathbf{u}, \mathbf{v}) \psi_l(\mathbf{v}) d\mathbf{v},$$

where $\mathbf{u} = \mathbf{u}(t)$ and $\mathbf{v} = \mathbf{v}(t)$, and an RKHS can be formally defined as the closure of a linear combination of basis functions $\{\sqrt{\lambda_l} \psi_l(\mathbf{v})\}_{l=1}^\infty$ (Pillai et al. 2007). One may conduct inference in an RKHS by assuming a GP prior distribution over the functional covariates directly (Rasmussen and Williams 2006),

$$f(\mathbf{F}_i(t)) \sim \mathcal{GP}(\mu(\mathbf{F}_i(t)), \sigma(\mathbf{F}_i(t), \mathbf{F}_j(t))), \quad i, j = 1, \dots, n, \quad (8)$$

where $f(\bullet)$ is a smooth operator from \mathbb{L}^2 to \mathbb{R} that is completely specified by its mean function and positive definite covariance function, $\mu(\bullet)$ and $\sigma(\bullet, \bullet)$, respectively. Recall from Section 2 that, in practical applications, there are a finite number of observed topological summary statistics taken from a given geometric object. Therefore, if we condition on these finite set of locations, the prior distribution in (8) may be represented as multivariate normal (Kolmogorov and Rozanov 1960). Consider the following joint “weight-space” probabilistic regression model to complete our specification (Rasmussen and Williams 2006),

$$\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon}, \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \Sigma(\mathbf{F}(t), \mathbf{F}(t))), \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}), \quad (9)$$

where $\mathbf{f} = [f(\mathbf{F}_1(t)), \dots, f(\mathbf{F}_n(t))]^\top$ is now assumed to come from a multivariate normal with mean $\mathbf{0}$ (for simplicity) and covariance matrix $\Sigma(\mathbf{F}(t), \mathbf{F}(t))$.

3.2. Posterior Predictive Inference

We now formally describe how to conduct posterior predictive inferences on clinical phenotypic traits for unobserved patients. Assume that we have received a set of new brain tumors and have computed their corresponding topological summary statistics $\mathbf{F}^*(t)$. Under the prior in Equation (8), we can write the joint distribution between the observed patient responses (\mathbf{y}) and the function values taken at the test images (\mathbf{f}^*) as

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}^* \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \Sigma(\mathbf{F}(t), \mathbf{F}(t)) + \tau^2 \mathbf{I} & \Sigma(\mathbf{F}(t), \mathbf{F}^*(t)) \\ \Sigma(\mathbf{F}^*(t), \mathbf{F}(t)) & \Sigma(\mathbf{F}^*(t), \mathbf{F}^*(t)) \end{bmatrix}\right). \quad (10)$$

Intuitively, if we train the model on n tumors and there are n^* test images, then $\Sigma(\mathbf{F}(t), \mathbf{F}^*(t))$ results in an $n \times n^*$ matrix of covariances between each of the training and testing points. Similar interpretations also hold for other covariance entries as well. Deriving the conditional distributions for Equation (10) then results in a multivariate normal posterior predictive distribution for the test shape smooth operators $\mathbf{f}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$,

where

$$\begin{aligned} \boldsymbol{\mu}^* &= \Sigma(\mathbf{F}^*(t), \mathbf{F}(t)) [\Sigma(\mathbf{F}(t), \mathbf{F}(t)) + \tau^2 \mathbf{I}]^{-1} \mathbf{y} \\ \boldsymbol{\Sigma}^* &= \Sigma(\mathbf{F}^*(t), \mathbf{F}^*(t)) - \Sigma(\mathbf{F}^*(t), \\ &\quad \mathbf{F}(t)) [\Sigma(\mathbf{F}(t), \mathbf{F}(t)) + \tau^2 \mathbf{I}]^{-1} \Sigma(\mathbf{F}(t), \mathbf{F}^*(t)). \end{aligned} \tag{11}$$

Note that in many applications, covariance functions can be indexed by a bandwidth or length-scale parameter θ , $\sigma_\theta(\mathbf{u}, \mathbf{v})$. For example, the Gaussian kernel can be specified as $\sigma_\theta(\mathbf{u}, \mathbf{v}) = \exp\{-\|\mathbf{u} - \mathbf{v}\|^2/2\theta\}$. This bandwidth parameter can be inferred; however, posterior inference over θ is slow, complicated, and often mixes poorly (Liang et al. 2009). For simplicity, we will work with a fixed bandwidth that is chosen via 10-fold cross-validation.

4. Predicting Clinical Outcomes in Glioblastoma

To fully illustrate the statistical utility of tumor images (captured by the SECT topological summary statistic), we apply the GP regression model to a GBM radiomic study with two measured clinical outcomes: disease free survival (DFS) and overall survival (OS). Some recent work in radiomics has confirmed the utility of imaging data in GBM research. These efforts suggest that the inclusion of shape information improves both the prediction of patient survival outcomes, as well as the classification of tumor subtypes (Gutman et al. 2013; Mazurowski, Desjardins, and Malof 2013; Gevaert et al. 2014; Macyszyn et al. 2016). It is important to distinguish, however, that most of these previous studies were limited to gross spatial features of cancer images (e.g., the presence of multifocal tumors, the location of recurrent lesions, or crude volumetric calculations). The SECT offers a novel contribution to radiomic research as a topological representation of imaging data. In this section, we will specifically assess whether topological features are better predictors of DFS and OS prognoses than three other tumor characteristics: (i) gene expression, (ii) tumor morphometry, and (iii) tumor geometry.

4.1. Genomic and Radiomic Data

MRIs of primary GBM tumors were collected from $n = 48$ patients archived by The Cancer Imaging Archive (TCIA) (Clark et al. 2013; Scarpace et al. 2016), which is a publicly accessible data repository containing medical images of cancer patients with matched genomic and clinical data collected by The Cancer Genome Atlas (TCGA) (The Cancer Genome Atlas Research Network 2008). These patients were selected based on two criteria: (i) individuals had post-contrast T1 axial (transverse) MRIs taken at the time of their diagnosis, and (ii) these patients have matching (mRNA) gene expression data and clinical correlates (e.g., recorded DFS and OS) that are publicly available (Gao et al. 2013; Grossman et al. 2016). There are three key factors that influenced our decision to use this particular subset of samples. First, the T1-weighted MRI with gadolinium contrast is one of the most commonly used imaging modalities and is often implemented to assay lesions with vascular activity (Adin et al. 2015). Second, axial (transverse) slices were considered as this was the most common representation in TCIA database and resulted in the

dataset with the most observations. Third, exclusively using MRIs taken at the time of diagnosis allowed us to avoid any potential confounding factors related to treatment-specific effects that may alter postoperative imaging and/or genomic profiles (Macyszyn et al. 2016).

Each collection of patient MRIs consisted of approximately 23–25 segmented slices of two-dimensional grayscale images (with the exact numbers varying between patients). We segment these images with the computer-assisted program MITKats to extract tumor lesions from the surrounding brain tissue (Chen and Rabadán 2017). Briefly, this algorithm first converts MRI images to a grayscale, and then thresholds to generate binary images. Morphological segmentation is then applied to delineate connected components. This is done by selecting contours corresponding to enhanced tumor lesions, which are lighter than healthy brain tissue. For instance, necrosis (or H_1 or H_2 homology as described above in Section 2.2) is represented by dark regions nested within the indicated lesion. An example of a raw image obtained from TCIA, along with its corresponding segmentation, is given in Figure S4.

From these segmented images, we collect three types of tumor shape information: morphometric features, geometric measurements, and topological summary statistics. Here, we use the same morphometric features outlined in previous imaging studies (Han et al. 2010; Chang et al. 2011) (see listed references for specific details on extraction and computation). The final dataset consisted of these 212 morphometric predictors corresponding to shape and texture, including cellularity skewness, cytoplasm intensity, nucleus texture, nucleus curvature, and median edge length. We also consider five tumor geometric measures. The first is the enhancing volume for each MRI slice, which is summed over lesions in the multifocal case. The other geometric measurements are the core volume of the enhancing and necrotic regions, the longest lesion diameter, and the shape factor of the tumor. For the purposes of this study, we define the shape factor to be the longest lesion diameter divided by the diameter of a sphere with the same volume. Lastly, when computing topological summary statistics, we use 100 sublevel sets per slice and compute a different EC curve for 72 directions evenly sampled over the interval $[0, 2\pi]$. After concatenation, this results in 7200-dimensional vectors for each patient. Averaging over these curves for each direction gives the proposed smooth EC statistics. It is well known that reconstructing three-dimensional brain tissue (and corresponding tumors) from two-dimensional slices is a nontrivial task (Cline et al. 1987; Amruta, Gole, and Karunakar 2010; Jaffar et al. 2012). Moreover, in the context of our case study, it is not guaranteed that the space in-between individual slices will be the same for each patient. Therefore, we aggregate topological summaries across slices for each direction as a proxy for rotating (accurate) three-dimensional representations of each tumor. Example SECT summary statistics for a segmented tumor are depicted in Figure 3.

Finally, we use matched mRNA gene expression levels of the preselected TCGA samples as a baseline data source. Following specific preprocessing steps from other genomic studies (Singleton et al. 2017), we use the robust multiarray average (RMA) normalization procedure to correct for potential lab-based batch effects and other potential confounders (Irizarry

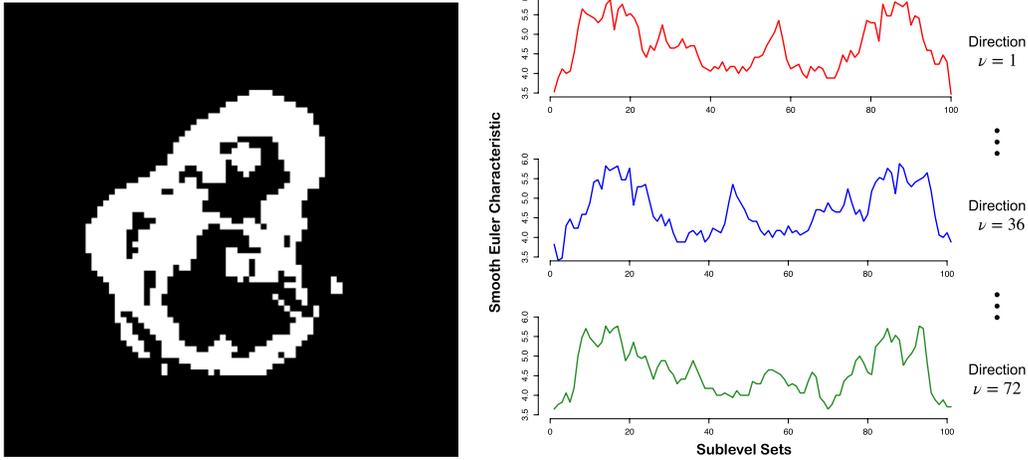


Figure 3. Example of a segmented tumor (on the left) with its corresponding SECT curves (on the right). On the right side, we illustrate the curves for directions $\nu = 1$, 36, and 72, respectively. The x-axis shows the number of sublevel sets (i.e., steps) as the filtration progresses, which we fix to be 100 per slice. In the radiomic prediction analysis, all 72 curves are concatenated together to create a 7200-dimensional covariate vector for the patient.

et al. 2003). This resulted in a final dataset consisting of 8725 genes, which also passed a prespecified hybridization accuracy threshold and showed reasonably varying expression across the assay.

4.2. Prediction Results

We now compare the ability of each data type to predict two clinical outcomes: DFS and OS. Briefly, DFS is the period after a successful treatment during which there are no signs or symptoms of the cancer, while OS is tabulated as the entire period after the initial treatment where the patient is still alive. It is worth noting that DFS is more commonly used over OS in adjuvant cancer clinical trials because it offers earlier presentation of data (Sargent et al. 2005). This stems from the idea that events due to disease recurrence occur earlier than death from disease, thus resulting in a cleaner signal (Birgisson et al. 2011).

Each tumor feature type is modeled with the GP regression model detailed in Section 3. In the context of this application, every patient in the data has an official death time. Hence, there is no need for a right-censored analysis or Cox-based methods (Rutledge et al. 2013; Fatai and Gamielien 2018; Kundu et al. 2018). We use two types of metrics to assess prediction accuracy: (i) the squared correlation coefficient (R^2) and (ii) the frequency for which a given data type exhibits the greatest R^2 , which we denote as Optimal%. When analyzing each outcome, we randomly split the data 1000 different times into 80% training and 20% out-of-sample test sets. In each case, the survival times are centered and scaled to have mean 0 and variance 1 to facilitate the interpretation of results. To illustrate the robustness of the SECT, we apply the GP regression framework using three different covariance functions. The goal is to show that the power of the SECT summary statistic is robust to this choice. Here, suppose that \mathbf{u} and \mathbf{v} are two different covariate vectors. We consider the linear (gram) kernel $\sigma(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v} / p$, where p is the length of the feature vector for a given data type; the Gaussian kernel $\sigma_\theta(\mathbf{u}, \mathbf{v}) = \exp\{-\|\mathbf{u} - \mathbf{v}\|^2 / 2\theta\}$; and the Cauchy

kernel $\sigma_\theta(\mathbf{u}, \mathbf{v}) = (1 + \theta \|\mathbf{u} - \mathbf{v}\|^2)^{-1}$. As previously mentioned, the last two functions are indexed by a bandwidth or length-scale parameter θ , which we select via 10-fold cross-validation over the grid $[0.1, 10]$ with step sizes equal to 0.1. Briefly, a value of 0 denotes a rigid function, while 10 represents a smoothed estimator. In Table 1, we present the mean R^2 and corresponding standard errors across testing splits to show how each tumor characteristic performs while taking into account variability. This table also lists the estimated bandwidths that generated these results.

Overall, our study shows that SECT topological summaries result in the most accurate predictions for survival—particularly for DFS. For example, using the Cauchy kernel function with SECT features resulted in the greatest R^2 for both DFS and OS at 0.237 and 0.158, respectively. This led to the SECT being the optimal tumor characteristic 40.5% of the time for DFS and 36.5% of the time for OS. There are a few possible explanations for these results. First and foremost, gene expression is known to be highly variable, particularly in GBM (Verhaak et al. 2010). Second, volumetric-based measurements only detail information about tumor size, but this information is likely not enough to be an effective predictor of patient survival on its own. Instead, one could imagine the geometry of a tumor being more useful when paired with the spatial location of lesions; hence, better detailing the severity of the shape. Conversely, morphometric features describe more focal-based characteristics of malignancies. This attention to texture causes more global information about the tumor to be missed.

These predictive results may also be due to the nature of the clinical outcomes that we chose to model. As previously mentioned, DFS is a prognostic measure of cancer recurrence and corresponds to the reappearance of the disease after initial treatment. This correlate can often be better defined than OS, where the cause of a patient's death may not necessarily be due to cancer-based complications. Indeed, each of the tumor characteristics that we consider generally perform better when predicting DFS versus OS (see Table 1). Nonetheless, measurements that provided detailed information about one particular aspect

Table 1. Detailed results for predicting disease free survival (DFS) and overall survival (OS) using Gaussian process regression models defined by the linear, Gaussian, and Cauchy covariance functions, respectively.

Covariance function(s)	Data type	Disease free survival (DFS)			Overall survival (OS)		
		R^2	Optimal%	$\hat{\theta}$	R^2	Optimal%	$\hat{\theta}$
Linear kernel	Gene expression	0.097 (0.013)	18.1%	–	0.075 (0.01)	18.4%	–
	Morphometrics	0.133 (0.015)	23.6%	–	0.127 (0.015)	33.9%	–
	Geometrics	0.137 (0.017)	23.9%	–	0.100 (0.012)	23.2%	–
	SECT	0.198 (0.023)	34.4%	–	0.097 (0.012)	24.5%	–
Gaussian kernel	Gene expression	0.129 (0.015)	23.6%	4.3	0.082 (0.011)	16.2%	10
	Morphometrics	0.113 (0.013)	16.3%	0.1	0.113 (0.012)	22.3%	4.0
	Geometrics	0.154 (0.018)	21.1%	5.2	0.102 (0.013)	22.3%	5.0
	SECT	0.228 (0.026)	36.1%	0.6	0.168 (0.018)	39.2%	4.2
Cauchy kernel	Gene expression	0.126 (0.015)	26.2%	6.4	0.084 (0.010)	16.7%	10.0
	Morphometrics	0.088 (0.012)	15.9%	1.2	0.115 (0.014)	27.9%	4.5
	Geometrics	0.116 (0.013)	17.4%	0.2	0.095 (0.012)	18.9%	3.5
	SECT	0.237 (0.027)	40.5%	0.6	0.158 (0.017)	36.5%	5.5

NOTE: For each model fit, we consider the predictive utility of four different genomic data types: gene expression, tumor morphometry, tumor geometry, and the proposed smooth Euler characteristic transform (SECT). Assessment is carried out by using the predictive squared correlation coefficient (R^2), where larger numbers indicate better performance. We also use Optimal% to denote the percentage of the time that a model exhibits the greatest R^2 . All values in bold represent the best method in these two assessment categories. These values are based on 1000 random 80–20 splits for each clinical outcome. Standard errors for each model are given the parentheses. Lastly, we give estimates for the bandwidth or length-scale parameter θ used to compute each kernel function. Note that θ was found by using 10-fold cross-validation over the grid [0.1, 10] with step sizes equal to 0.1.

of the disease (e.g., volumetric and morphometric quantities) are not as relevant as those that aim to illustrate a more comprehensive view. Thus, topological features are effective predictors as they provide some notion about both the size and texture of tumors.

4.3. MRI Specific Consequences on Results

In our study, the robustness of the SECT to choice of metric is particularly relevant because the geometric structure of the brain is known to be fibrous—meaning that the brain is made up of, and connected by, cerebral fiber pathways (Wedeen et al. 2012). This brings into question the validity of assuming the usual Euclidean metric when quantifying shape. Both volumetric and morphometric analyses require the specification of a metric and, in the case where the usual assumption of a Euclidean measure does not apply, an appropriate one must be constructed. This is not always a straightforward task. Moreover, in fibrous settings, there is also the possibility for the further requirement of defining a geodesic. Examining topological properties, as opposed to metric-based properties, bypasses these technical difficulties. Altogether, incorporating a topological measure that is not based on a metric results in the flexibility to compare tumors of different sizes more seamlessly. Subsequently, this also implicitly allows for comparisons between different stages of the disease without needing to account for time of progression. We hypothesize that these flexible characteristics also contribute to the SECT being a better predictor of prognosis and survival.

4.4. More Biological Implications

One key implication from our results in DFS is that there possibly exist correlations between the topology of tumors and the molecular heterogeneity arising from the activation of dif-

ferent recurrence mechanisms. An example of this relationship occurs in (multifocal) tumors where lesions on opposite hemispheres of the brain originate from the same oncogenic effects, but events such as therapeutic resistance or cancer recurrence happen in only one hemisphere. This variation can be clinically relevant. It was recently proposed that these type of topologically based traits are linked to the mutation status of certain oncogenic relapse drivers (Lee et al. 2017). Hence, there is growing evidence that potential pathways of progression in GBM should go beyond the simple consideration of physical proximity (i.e., closeness in a geometric sense). For instance, a particular path to recurrence in GBM may be due to ambient effects inherent to a particular hemisphere of the brain. The prediction results we present in this work suggest that the topological features extracted by the SECT may be better than simple geometric summaries at providing insight into biological phenomena at the molecular level.

5. Discussion

In this article, we sought to quantify images of GBM tumors given by MRIs for statistical analyses and to demonstrate the clinical relevance of this information. To this end, we developed a topological summary statistic transform which maps shapes into a space that admits an inner product structure that is amenable to standard functional and nonlinear regression models. We then used our summary transform to predict the survival of GBM patients using our specified functional GP regression model. In this study, we compared the predictive accuracy using both molecular biomarkers and shape covariates. The SECT was shown to explain more of the variance in DFS of patients than all other covariates in a wide variety of models defined by various kernel functions. For the Gaussian and Cauchy kernels, in particular, the SECT outperformed the other measures in accounting for the variance in both DFS and OS.

Despite these results, several interesting future directions and open questions still remain. For example, in the current study, we focus solely on measuring how well topological features predict survival. Many studies in the radiomics space use deep learning approaches for accurate classification and prediction-based tasks (Lao et al. 2017; Li et al. 2017; Bibault et al. 2018; Rathore et al. 2018). Unfortunately, in this work, we did not have access to data with large enough sample sizes for the effective training of neural networks. However, in the future, it would be useful to see how our topological summary statistics may be integrated within deep learning frameworks. To ensure power, utilizing protected data from current consortium studies with a large number of participants would be of high interest (Mueller et al. 2005; Gounder et al. 2015; Sudlow et al. 2015). Moving away from prediction, it would also be useful to infer which particular spatial regions of the tumor are most relevant to clinical outcomes. Recent variable selection approaches for kernel-based methods can be used to infer the directions and segments of the Euler curves that are most relevant (Crawford et al. 2018, 2019). In this case, an important open problem is having the ability to recover, or partially reconstruct, a shape based on significant SECT summary statistics. Similarly, the distance measure for the SECT stated in Equation (7) provides a framework for comparing the shapes of tumors, and correlating geometric properties with molecular and clinical features. Understanding the relationship between therapeutic strategies, signaling pathway dependence, and tumor shapes would provide useful information about different forms of GBM and their etiologies. We conjecture that greater general knowledge about tumor shape may help in distinguishing true progression from pseudoprogression. Here, progression refers to the growth of the tumor itself, while a pseudoprogressing tumor has been infiltrated by immune cells and other factors.

Data Availability

The results shown here are in whole or part based upon data generated by the TCGA Research Network (<http://cancergenome.nih.gov/>). DICOM formatted MRI scans and patient clinical information were taken directly from the TCIA web portal (<https://wiki.cancerimagingarchive.net/display/Public/TCGA-GBM>). Matched molecular data were downloaded directly from the Genomic Data Commons (GDC) by selecting the RNA-Seq tab option (<https://portal.gdc.cancer.gov/projects/TCGA-GBM>). Shape-based summary statistics necessary for replicating this study (i.e., the segmented tumor images, the volumetric measurements, morphometric data, and topological summary statistics) are also publicly available on the SECT GitHub repository.

Software Availability

Software to compute the SECT from images and fit the GP regression model is publicly available in both R and MATLAB code, and located on the repository <https://github.com/lorinanthony/SECT>. The MRI images were segmented using the Medical Imaging Interaction Toolkit with augmented tools for segmentation (MITKats), which was written C++ and is

located at <https://github.com/RabadanLab/MITKats> (Chen and Rabadán 2017).

Supplementary Materials

The supplementary materials for this article consist of supplementary images to illustrate various mathematical concepts described in the text, as well as a more detailed discussion on the mathematics underlying the construction of the SECT.

Acknowledgments

The authors wish to thank Mao Li (Donald Danforth Plant Science Center) and Christoph Hellmayr (Duke University) for help with the formulation of code, as well as Francesco Abate (McKinsey & Co.), Katharine Turner (Australian National University), and Jiguang Wang (Hong Kong University of Science and Technology) for helpful conversations and input on a previous version of the manuscript. The authors would also like to acknowledge The Cancer Imaging Archive (TCIA) and The Cancer Genome Atlas (TCGA) initiatives for making the imaging and the clinical data used in this study publicly available.

Disclosure Statement

The authors have declared that no competing interests exist.

Funding

During some of this work, LC, AM, and RR were supported by the National Cancer Institute Physical Sciences–Oncology Network (NCI PS–ON) under grant no. 5 U54 CA 193313-02. AM was the PI on Pilot grant subaward no. G11124 for research on radiomics and radiogenomics. AM is also supported by the Irving Institute's CaMPR initiative under grant no. GG011557, and would like to acknowledge the support of the New Frontiers in Research Fund–Fonds Nouvelles Frontières en Recherche (SSHRC–NFRF–FNFR Government of Canada) NFRFE-2018-00431. LC would like to acknowledge the support of grants P20GM109035 (COBRE Center for Computational Biology of Human Disease; PI Rand) and P20GM103645 (COBRE Center for Central Nervous; PI Sanes) from the NIH NIGMS, 2U10CA180794-06 from the NIH NCI and the Dana Farber Cancer Institute (PIs Gray and Gatsonis), and an Alfred P. Sloan Research Fellowship (no. FG-2019-11622). AXC would like to acknowledge support by the Columbia University Medical Scientist Training Program (MSTP). SM would like to acknowledge funding from NSF DEB-1840223, NIH R01 DK116187-01, HFSP RGP0051/2017, NSF DMS 17-13012, and NSF DMS 16-13261. This work used a high-performance computing facility partially supported by grant 2016-IDG-1013 (“HARDAC+: Reproducible HPC for Next-generation Genomics”) from the North Carolina Biotechnology Center. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of any of the funders.

References

- Adin, M. E., Kleinberg, L., Vaidya, D., Zan, E., Mirbagheri, S., and Yousef, D. M. (2015), “Hyperintense Dentate Nuclei on T1-Weighted MRI: Relation to Repeat Gadolinium Administration,” *American Journal of Neuroradiology*, 36, 1859–1865. [1145]
- Al-Aifari, R., Daubechies, I., and Lipman, Y. (2013), “Continuous Procrustes Distance Between Two Surfaces,” *Communications on Pure and Applied Mathematics*, 66, 934–964. [1140]
- Amruta, A., Gole, A., and Karunakar, Y. (2010), “A Systematic Algorithm for 3-D Reconstruction of MRI Based Brain Tumors Using Morphological Operators and Bicubic Interpolation,” in *2010 2nd International Conference on Computer Technology and Development*, pp. 305–309. [1145]

- Barker, F. G., Davis, R. L., Chang, S. M., and Prados, M. D. (1996), "Necrosis as a Prognostic Factor in Glioblastoma Multiforme," *Cancer*, 77, 1161–1166. [1140]
- Bibault, J.-E., Giraud, P., Durdux, C., Taieb, J., Berger, A., Coriat, R., Chausade, S., Dousset, B., Nordlinger, B., and Burgun, A. (2018), "Deep Learning and Radiomics Predict Complete Response After Neo-Adjuvant Chemoradiation for Locally Advanced Rectal Cancer," *Scientific Reports*, 8, 12611. [1148]
- Birgisson, H., Wallin, U., Holmberg, L., and Glimelius, B. (2011), "Survival Endpoints in Colorectal Cancer and the Effect of Second Primary Other Cancer on Disease Free Survival," *BMC Cancer*, 11, 438. [1146]
- Bookstein, F. L. (1997), *Morphometric Tools for Landmark Data: Geometry and Biology*, Cambridge: Cambridge University Press. [1140]
- Boyer, D. M., Lipman, Y., St. Clair, E., Puente, J., Patel, B. A., Funkhouser, T., Jernvall, J., and Daubechies, I. (2011), "Algorithms to Automatically Quantify the Geometric Similarity of Anatomical Surfaces," *Proceedings of the National Academy of Sciences of the United States of America*, 108, 18221–18226. [1140]
- Boyer, D. M., Puente, J., Gladman, J. T., Glynn, C., Mukherjee, S., Yapunchich, G. S., and Daubechies, I. (2015), "A New Fully Automated Approach for Aligning and Comparing Shapes," *The Anatomical Record*, 298, 249–276. [1140]
- Chang, H., Fontenay, G. V., Han, J., Cong, G., Baehner, F. L., Gray, J. W., Spellman, P. T., and Parvin, B. (2011), "Morphometric Analysis of TCGA Glioblastoma Multiforme," *BMC Bioinformatics*, 12, 484. [1145]
- Chen, A. X., and Rabadán, R. (2017), "A Fast Semi-Automatic Segmentation Tool for Processing Brain Tumor Images," in *Towards Integrative Machine Learning and Knowledge Extraction*, eds. A. Holzinger, R. Goebel, M. Ferri, and V. Palade, Cham: Springer International Publishing, pp. 170–181. [1145,1148]
- Clark, K., Vendt, B., Smith, K., Freymann, J., Kirby, J., Koppel, P., Moore, S., Phillips, S., Maffitt, D., Pringle, M., Tarbox, L., and Prior, F. (2013), "The Cancer Imaging Archive (TCIA): Maintaining and Operating a Public Information Repository," *Journal of Digital Imaging*, 26, 1045–1057. [1139,1145]
- Cline, H. E., Dumoulin, C. L., Hart, H. R. J., Lorensen, W. E., and Ludke, S. (1987), "3D Reconstruction of the Brain From Magnetic Resonance Images Using a Connectivity Algorithm," *Magnetic Resonance Imaging*, 5, 345–352. [1145]
- Crawford, L., Flaxman, S., Runcie, D., and West, M. (2019), "Variable Prioritization in Nonlinear Black Box Methods: A Genetic Association Case Study," *The Annals of Applied Statistics*, 13, 958–989. [1148]
- Crawford, L., Wood, K. C., Zhou, X., and Mukherjee, S. (2018), "Bayesian Approximate Kernel Regression With Variable Selection," *Journal of the American Statistical Association*, 113, 1710–1721. [1148]
- Dryden, I., and Mardia, K. (1998), *Statistical Shape Analysis*, Wiley Series in Probability and Statistics, New York: Wiley. [1140]
- Edelsbrunner, H., and Harer, J. (2010), *Computational Topology: An Introduction*, Providence, RI: American Mathematical Society. [1142]
- Edelsbrunner, H., Letscher, D., and Zomorodian, A. (2000), "Topological Persistence and Simplification," in *41st Annual Symposium on Foundations of Computer Science, FOCS'00*, IEEE Computer Society, Washington, DC, USA. [1140]
- Fatai, A. A., and Gamielidien, J. (2018), "A 35-Gene Signature Discriminates Between Rapidly- and Slowly-Progressing Glioblastoma Multiforme and Predicts Survival in Known Subtypes of the Cancer," *BMC Cancer*, 18, 377. [1146]
- Ferraty, F., and Vieu, P. (2006), *Nonparametric Functional Data Analysis: Theory and Practice*, New York: Springer Science & Business Media. [1144]
- Friston, K., Phillips, J., Chawla, D., and Buchel, C. (2000), "Nonlinear PCA: Characterizing Interactions Between Modes of Brain Activity," *Philosophical Transactions of the Royal Society of London, Series B*, 355, 135–146. [1144]
- Gao, J., Aksoy, B. A., Dogrusoz, U., Dresdner, G., Gross, B., Sumer, S. O., Sun, Y., Jacobsen, A., Sinha, R., Larsson, E., Cerami, E., Sander, C., and Schultz, N. (2013), "Integrative Analysis of Complex Cancer Genomics and Clinical Profiles Using the cBioPortal," *Science Signaling*, 6, pl1. [1145]
- Gao, T., Yapunchich, G. S., Daubechies, I., Mukherjee, S., and Boyer, D. M. (2018), "Development and Assessment of Fully Automated and Globally Transitive Geometric Morphometric Methods, With Application to a Biological Comparative Dataset With High Interspecific Variation," *The Anatomical Record*, 301, 636–658. [1140]
- Gevaert, O., Mitchell, L. A., Achrol, A. S., Xu, J., Echegaray, S., Steinberg, G. K., Cheshier, S. H., Napel, S., Zaharchuk, G., and Plevritis, S. K. (2014), "Glioblastoma Multiforme: Exploratory Radiogenomic Analysis by Using Quantitative Image Features," *Radiology*, 273, 168–174. [1145]
- Gounder, M. M., Nayak, L., Sahebjam, S., Muzikansky, A., Sanchez, A. J., Desideri, S., Ye, X., Ivy, S. P., Nabors, L. B., Prados, M., Grossman, S., DeAngelis, L. M., and Wen, P. Y. (2015), "Evaluation of the Safety and Benefit of Phase I Oncology Trials for Patients With Primary CNS Tumors," *Journal of Clinical Oncology*, 33, 3186–3192. [1148]
- Grossman, R. L., Heath, A. P., Ferretti, V., Varmus, H. E., Lowy, D. R., Kibbe, W. A., and Staudt, L. M. (2016), "Toward a Shared Vision for Cancer Genomic Data," *New England Journal of Medicine*, 375, 1109–1112. [1145]
- Gutman, D. A., Cooper, L. A. D., Hwang, S. N., Holder, C. A., Gao, J., Aurora, T. D., Dunn, W. D., Scarpace, L., Mikkelsen, T., Jain, R., Wintermark, M., Jilwan, M., Raghavan, P., Huang, E., Clifford, R. J., Mongkolwat, P., Kleper, V., Freymann, J., Kirby, J., Zinn, P. O., Moreno, C. S., Jaffe, C., Colen, R., Rubin, D. L., Saltz, J., Flanders, A., and Brat, D. J. (2013), "MR Imaging Predictors of Molecular Profile and Survival: Multi-Institutional Study of the TCGA Glioblastoma Data Set," *Radiology*, 267, 560–569. [1145]
- Han, J., Chang, H., Andarawewa, K., Yaswen, P., Barcellos-Hoff, M. H., and Parvin, B. (2010), "Multidimensional Profiling of Cell Surface Proteins and Nuclear Markers," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 7, 80–90. [1145]
- Irizarry, R., Hobbs, C., Beazer-Barclay, Y., Antonellis, K., Scherf, U., and Speed, T. (2003), "Exploration, Normalization, and Summaries of High Density Oligonucleotide Array Probe Level Data," *Biostatistics*, 2, 249–264. [1146]
- Jaffar, M. A., Zia, S., Latif, G., Mirza, A. M., Mehmood, I., Ejaz, N., and Baik, S. W. (2012), "Anisotropic Diffusion Based Brain MRI Segmentation and 3D Reconstruction," *International Journal of Computational Intelligence Systems*, 5, 494–504. [1145]
- Kadri, H., Duflos, E., Preux, P., Canu, S., and Davy, M. (2010), "Nonlinear Functional Regression: A Functional RKHS Approach," in *Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS'10)* (Vol. 9), pp. 374–380. [1144]
- Kendall, D. G. (1984), "Shape Manifolds, Procrustean Metrics, and Complex Projective Spaces," *Bulletin of the London Mathematical Society*, 16, 81–121. [1140]
- Kolmogorov, A. N., and Rozanov, Y. A. (1960), "On Strong Mixing Conditions for Stationary Gaussian Processes," *Theory of Probability & Its Applications*, 5, 204–208. [1144]
- Kundu, S., Cheng, Y., Shin, M., Manyam, G., Mallick, B. K., and Baladandayuthapani, V. (2018), "Bayesian Variable Selection With Graphical Structure Learning: Applications in Integrative Genomics," *PLoS ONE*, 13, e0195070. [1146]
- Kusano, G., Fukumizu, K., and Hiraoka, Y. (2018), "Kernel Method for Persistence Diagrams via Kernel Embedding and Weight Factor," *The Journal of Machine Learning Research*, 18, 1–41. [1142]
- Kwitt, R., Huber, S., Niethammer, M., Lin, W., and Bauer, U. (2015), "Statistical Topological Data Analysis—A Kernel Perspective," *Advances in Neural Information Processing Systems*, 28, 3070–3078. [1142]
- Lao, J., Chen, Y., Li, Z.-C., Li, Q., Zhang, J., Liu, J., and Zhai, G. (2017), "A Deep Learning-Based Radiomics Model for Prediction of Survival in Glioblastoma Multiforme," *Scientific Reports*, 7, 10353. [1148]
- Lee, J.-K., Wang, J., Sa, J. K., Ladewig, E., Lee, H.-O., Lee, I.-H., Kang, H. J., Rosenbloom, D. S., Camara, P. G., Liu, Z., van Nieuwenhuizen, P., Jung, S. W., Choi, S. W., Kim, J., Chen, A., Kim, K.-T., Shin, S., Seo, Y. J., Oh, J.-M., Shin, Y. J., Park, C.-K., Kong, D.-S., Seol, H. J., Blumberg, A., Lee, J.-I., Iavarone, A., Park, W.-Y., Rabadan, R., and Nam, D.-H. (2017), "Spatiotemporal Genomic Architecture Informs Precision Oncology in Glioblastoma," *Nature Genetics*, 49, 594–599. [1147]
- Li, Z., Wang, Y., Yu, J., Guo, Y., and Cao, W. (2017), "Deep Learning Based Radiomics (DLR) and Its Usage in Noninvasive IDH1 Prediction for Low Grade Glioma," *Scientific Reports*, 7, 5467. [1148]

- Liang, F., Mao, K., Mukherjee, S., and West, M. (2009), "Nonparametric Bayesian Kernel Models," Technical Report, Duke University, Department of Statistical Science. [1145]
- Lipman, Y., and Daubechies, I. (2011), "Conformal Wasserstein Distances: Comparing Surfaces in Polynomial Time," *Advances in Mathematics*, 227, 1047–1077. [1140]
- Liu, Q., Liu, Y., Li, W., Wang, X., Sawaya, R., Lang, F. F., Yung, W. K., Chen, K., Fuller, G. N., and Zhang, W. (2015), "Genetic, Epigenetic, and Molecular Landscapes of Multifocal and Multicentric Glioblastoma," *Acta Neuropathologica*, 130, 587–597. [1140]
- Macyszyn, L., Akbari, H., Pisapia, J. M., Da, X., Attiah, M., Pigrish, V., Bi, Y., Pal, S., Davuluri, R. V., Rococgrandi, L., Dahmane, N., Martinez-Lage, M., Biros, G., Wolf, R. L., Bilello, M., O'Rourke, D. M., and Davatzikos, C. (2016), "Imaging Patterns Predict Patient Survival and Molecular Subtype in Glioblastoma via Machine Learning Techniques," *Neuro-Oncology*, 18, 417–425. [1145]
- Mazurowski, M. A., Desjardins, A., and Malof, J. M. (2013), "Imaging Descriptors Improve the Predictive Power of Survival Models for Glioblastoma Patients," *Neuro-Oncology*, 15, 1389–1394. [1145]
- Mercer, J. (1909), "Functions of Positive and Negative Type and Their Connection With the Theory of Integral Equations," *Philosophical Transactions of the Royal Society of London, Series A*, 209, 415–446. [1144]
- Mileyko, Y., Mukherjee, S., and Harer, J. (2011), "Probability Measures on the Space of Persistence Diagrams," *Inverse Problems*, 27, 124007. [1142]
- Morris, J. S. (2015), "Functional Regression," *Annual Review of Statistics and Its Application*, 2, 321–359. [1144]
- Mueller, S. G., Weiner, M. W., Thal, L. J., Petersen, R. C., Jack, C., Jagust, W., Trojanowski, J. Q., Toga, A. W., and Beckett, L. (2005), "The Alzheimer's Disease Neuroimaging Initiative," *Neuroimaging Clinics*, 15, 869–877. [1148]
- Müller, H.-G., and Stadtmüller, U. (2005), "Generalized Functional Linear Models," *The Annals of Statistics*, 33, 774–805. [1144]
- Müller, H.-G., and Yao, F. (2008), "Functional Additive Models," *Journal of the American Statistical Association*, 103, 1534–1544. [1144]
- Patel, M. A., Kim, J. E., Ruzevick, J., Li, G., and Lim, M. (2014), "The Future of Glioblastoma Therapy: Syngism of Standard of Care and Immunotherapy," *Cancers*, 6, 1953–1985. [1139]
- Pillai, N. S., Wu, Q., Liang, F., Mukherjee, S., and Wolpert, R. (2007), "Characterizing the Function Space for Bayesian Kernel Models," *Journal of Machine Learning Research*, 8, 1769–1797. [1144]
- Pomann, G.-M., Staicu, A.-M., Lobaton, E. J., Mejia, A. F., Dewey, B. E., Reich, D. S., Sweeney, E. M., and Shinohara, R. T. (2016), "A Lag Functional Linear Model for Prediction of Magnetization Transfer Ratio in Multiple Sclerosis Lesions," *The Annals of Applied Statistics*, 10, 2325–2348. [1144]
- Preda, C. (2007), "Regression Models for Functional Data by Reproducing Kernel Hilbert Spaces Methods," *Journal of Statistical Planning and Inference*, 137, 829–840. [1144]
- Ramsay, J. O. (2006), *Functional Data Analysis*, Wiley Online Library. [1144]
- Rasmussen, C. E., and Williams, C. K. I. (2006), *Gaussian Processes for Machine Learning*, Cambridge, MA: MIT Press. [1144]
- Rathore, S., Akbari, H., Rozycki, M., Abdullah, K. G., Nasrallah, M. P., Binder, Z. A., Davuluri, R. V., Lustig, R. A., Dahmane, N., Bilello, M., O'Rourke, D. M., and Davatzikos, C. (2018), "Radiomic MRI Signature Reveals Three Distinct Subtypes of Glioblastoma With Different Clinical and Molecular Characteristics, Offering Prognostic Value Beyond IDH1," *Scientific Reports*, 8, 5087. [1148]
- Reininghaus, J., Huber, S., Bauer, U., and Kwitt, R. (2015), "A Stable Multi-Scale Kernel for Topological Machine Learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4741–4748. [1142]
- Rutledge, W. C., Kong, J., Gao, J., Gutman, D. A., Cooper, L. A. D., Appin, C., Park, Y., Scarpace, L., Mikkelsen, T., Cohen, M. L., Aldape, K. D., McLendon, R. E., Lehman, N. L., Miller, C. R., Schniederjan, M. J., Brennan, C. W., Saltz, J. H., Moreno, C. S., and Brat, D. J. (2013), "Tumor-Infiltrating Lymphocytes in Glioblastoma Are Associated With Specific Genomic Alterations and Related to Transcriptional Class," *Clinical Cancer Research*, 19, 4951–4960. [1146]
- Sargent, D., Wieand, H., Haller, D., Gray, R., Benedetti, J., Buysse, M., Labianca, R., Seitz, J., O'Callaghan, C., Francini, G., Grothey, A., O'Connell, M., Catalano, P., Blanke, C., Kerr, D., Green, E., Wolmark, N., Andre, T., Goldberg, R., and De Gramont, A. (2005), Disease-Free Survival Versus Overall Survival As a Primary End Point For Adjuvant Colon Cancer Studies: Individual Patient Data From 20,898 Patients on 18 Randomized Trials," *Journal of Clinical Oncology*, 23, 8664–8670. [1146]
- Scarpace, L., Mikkelsen, T., Cha, S., Rao, S., Tekchandani, S., Gutman, D., Saltz, J., Eickson, B., Pedano, N., and Flanders, A. (2016), "Radiology Data From the Cancer Genome Atlas Glioblastoma Multiforme [TCGA-GBM] Collection," *The Cancer Imaging Archive*, 11, 4. [1139,1145]
- Singleton, K. R., Crawford, L., Tsui, E., Manchester, H. E., Maertens, O., Liu, X., Libertini, M. V., Magpusao, A. N., Stein, E. M., Tingley, J. P., Frederick, D. T., Boland, G. M., Flaherty, K. T., McCall, S. J., Krepler, C., Sproesser, K., Herlyn, M., Adams, D. J., Locasale, J. W., Cichowski, K., Mukherjee, S., and Wood, K. C. (2017), "Melanoma Therapeutic Strategies That Select Against Resistance by Exploiting MYC-Driven Evolutionary Convergence," *Cell Reports*, 21, 2796–2812. [1145]
- Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., Liu, B., Matthews, P., Ong, G., Pell, J., Silman, A., Young, A., Sprosen, T., Peakman, T., and Collins, R. (2015), "UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age," *PLoS Medicine*, 12, e1001779. [1148]
- The Cancer Genome Atlas Research Network (2008), "Comprehensive Genomic Characterization Defines Human Glioblastoma Genes and Core Pathways," *Nature*, 455, 1061–1068. [1140,1145]
- Turner, K., Mileyko, Y., Mukherjee, S., and Harer, J. (2014), "Fréchet Means for Distributions of Persistence Diagrams," *Discrete & Computational Geometry*, 52, 44–70. [1140]
- Turner, K., Mukherjee, S., and Boyer, D. M. (2014), "Persistent Homology Transform for Modeling Shapes and Surfaces," *Information and Inference*, 3, 310–344. [1140,1141,1142,1143]
- Verhaak, R. G. W., Hoadley, K. A., Purdom, E., Wang, V., Qi, Y., Wilkerson, M. D., Miller, C. R., Ding, L., Golub, T., Mesirov, J. P., Alexe, G., Lawrence, M., O'Kelly, M., Tamayo, P., Weir, B. A., Gabriel, S., Winckler, W., Gupta, S., Jakkula, L., Feiler, H. S., Hodgson, J. G., James, C. D., Sarkaria, J. N., Brennan, C., Kahn, A., Spellman, P. T., Wilson, R. K., Speed, T. P., Gray, J. W., Meyerson, M., Getz, G., Perou, C. M., and Hayes, D. N. (2010), "Integrated Genomic Analysis Identifies Clinically Relevant Subtypes of Glioblastoma Characterized by Abnormalities in PDGFRA, IDH1, EGFR, and NF1," *Cancer Cell*, 17, 98–110. [1146]
- Wahba, G. (1997), "Support Vector Machines, Reproducing Kernel Hilbert Spaces and the Randomized GACV," *Advances in Neural Information Processing Systems*, 6, 69–87. [1144]
- Wang, J.-L., Chiou, J.-M., and Müller, H.-G. (2016), "Functional Data Analysis," *Annual Review of Statistics and Its Application*, 3, 257–295. [1144]
- Wedeen, V. J., Rosene, D. L., Wang, R., Dai, G., Mortazavi, F., Hagmann, P., Kaas, J. H., and Tseng, W.-Y. I. (2012), "The Geometric Structure of the Brain Fiber Pathways," *Science*, 335, 1628–1634. [1147]
- Yuan, M., and Cai, T. T. (2010), "A Reproducing Kernel Hilbert Space Approach to Functional Linear Regression," *The Annals of Statistics*, 38, 3412–3444. [1144]
- Zomorodian, A., and Carlsson, G. (2005), "Computing Persistent Homology," *Discrete & Computational Geometry*, 33, 249–274. [1140]