

Second-order spatial analysis of shapes of tumor cell nuclei

Ye Jin Choi, Sebastian Kurtek

Department of Statistics, The Ohio State University

Simeng Zhu

James Cancer Center, The Ohio State University

Karthik Bharath

School of Mathematical Sciences, University of Nottingham

November 13, 2025

Abstract

Intra-tumor heterogeneity driving disease progression is characterized by distinct growth and spatial proliferation patterns of cells and their nuclei within tumor and non-tumor tissues. A widely accepted hypothesis is that these spatial patterns are correlated with morphology of the cells and their nuclei. Nevertheless, tools to quantify the correlation, with uncertainty, are scarce, and the state-of-the-art is based on low-dimensional numerical summaries of the shapes that are inadequate to fully encode shape information. To this end, we propose a marked point process framework to assess spatial correlation among shapes of planar closed curves, which represent cell or nuclei outlines. With shapes of curves as marks, the framework is based on a mark-weighted K function, a second-order spatial statistic that accounts for the marks' variation by using test functions that capture only the shapes of cells and their nuclei. We then develop local and global hypothesis tests for spatial dependence between the marks using the K function. The framework is brought to bear on the cell nuclei extracted from histopathology images of breast cancer, where we uncover distinct correlation patterns that are consistent with clinical expectations.

Keywords: Functional marked point process; Statistical shape analysis; Elastic metric; Histopathology imaging

1 Introduction

The tumor microenvironment (TME) comprises tumor cells, stromal (supportive) cells, and other components, together driving abnormal, rapidly growing structures [Hanahan and Coussens, 2012]. While both the tumor cell and its nucleus offer valuable information for tumor characterization, research has focused on studying the nucleus due to its more direct association with genetic alterations [Hanahan and Weinberg, 2011]. Further, extracting cell boundaries from histopathology images is often unreliable, owing to overlapping cells and poorly defined cytoplasmic edges. Tumor cell nuclei, unlike those of normal cells, frequently exhibit morphological abnormalities such as indentations, folds and fragmentation across various cancer types [Singh and Lele, 2022]. These nuclear irregularities are not only visually distinctive, but also serve as strong predictors of cancer progression and treatment outcomes [Nafe et al., 2005a, de Andrea et al., 2011].

Pathologists have long relied on nuclear morphology for diagnostic purposes. In the specific context of breast cancer, histopathological examination of cancer tissue provides critical prognostic information that guides clinical treatment decisions. Histologic grading serves as a standardized assessment tool whereby pathologists evaluate tumor specimens from biopsies or surgical excisions. The current Nottingham grading system incorporates three key morphological parameters: (i) the extent of tubule formation, (ii) the degree of nuclear pleomorphism, and (iii) mitotic activity count [Bloom and Richardson, 1957, Elston and Ellis, 1991]. However, the inherent subjectivity in manually assessing these morphological features has raised concerns regarding inter-observer reproducibility, prompting efforts to develop more objective, quantitative assessment methods [van Dooijeweert et al., 2022].

As a result, many recent studies have relied on quantitative approaches for evaluating morphological variation in tumor cell nuclei, derived from histopathology images, to pursue several scientific and clinical goals, including automatic differentiation of malignant and benign cells [Fischer, 2020], classification of tumor subtypes [Beck et al., 2011], prediction of cancer prognosis and patient survival [Murphy et al., 1990], discovery of associations with genetic mutations [Sali et al., 2024], and forecasting of patient response to chemotherapy [Kather et al., 2019]. Further, numerous studies have indicated that spatial interactions

among cells in the TME also help differentiate cancer grades [Barua et al., 2018, Tsujikawa et al., 2020, Wang et al., 2022]. However, these studies focused on cell locations only. In contrast, Sali et al. [2024] and Lu et al. [2020] emphasized that intra-tumor or local morphological diversity of tumor cells’ nuclei is a key cancer trait and a strong prognostic indicator. To capture their shape, Singh et al. [2010] represented the three-dimensional structure of stromal cells using spherical harmonics and explored principal modes of shape variation in relation to duct proximity. However, they only distinguished between ‘close’ or ‘far’ cells. Similarly, Nafe et al. [2005b] used various shape descriptors and size, along with spatial information, e.g., neighboring nuclei or inter-nuclear distances, for tumor grade classification. In both cases, representing the shape of cells or their nuclei using a finite number of spherical harmonics or subjectively chosen descriptors cannot fully capture their morphological complexity, resulting in potentially significant loss of information.

1.1 Motivation and contributions

This paper is motivated by the following questions.

- (i) Are the shapes of nuclei of tumor cells correlated with their spatial configurations within the tumor microenvironment?
- (ii) How does the morphological variation depend on the spatial proximity between normal or tumor cell nuclei?

Answers to the two questions along with a procedure to quantify the accompanying uncertainty offer crucial empirical validation to clinical observations. For example, in the context of breast cancer to be considered in Section 7, it is observed that normal breast cells arrange in a regular spatial tubular pattern around ducts while tumor cells form clusters based on similar genetic information Greaves and Maley [2012]; our findings using the developed framework corroborate these assertions.

Motivated by and tailored to tackle the two questions arising in the quantitative study of intra-tumor heterogeneity in two-dimensional pathology images, we propose a spatial point process model for the cell nuclei with marks that encode shape-related characteristics of their outlines. This is sensible since the location and number of cells (and their nuclei)

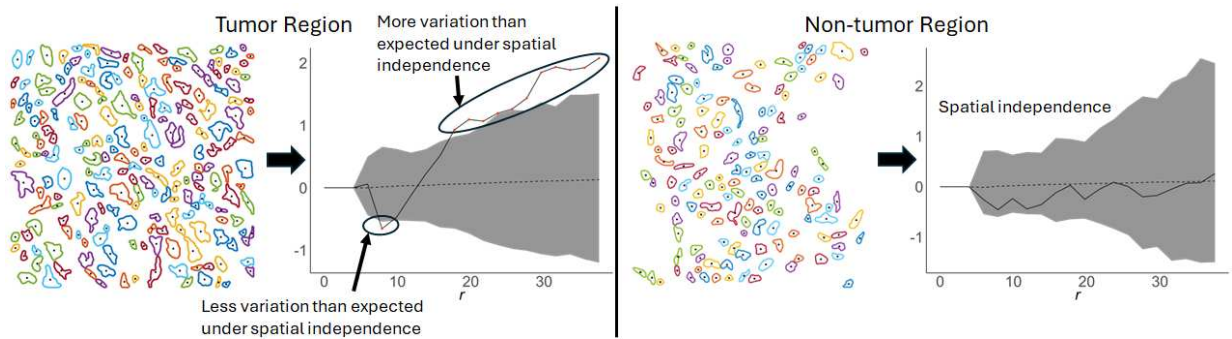


Figure 1: We contrast spatial correlation patterns at different scales among cell nuclei shapes within tumor (left) and non-tumor (right) regions via a second-order spatial statistic (solid black line). The shaded region is an envelope that captures underlying uncertainty.

in a histopathology image are not known *a priori*. The shape of a cell nucleus under the parameterized curve representation is a quantity unaffected by the nucleus' location, size, orientation and parameterization, referred to as its symmetries. While the physical location and parameterization of the nuclei constitute nuisance variation, size and orientation may provide useful information to characterize their morphological diversity, e.g., as related to tumor growth processes. We provide an illustration of the framework in Figure 1. We start with cell nuclei outlines within a region of interest. We then compute a second-order spatial statistic (black solid line), as described in Sections 4 and 5, that accounts for the cell nuclei' shape variation at different scales (defined through balls with different radii). In addition, we quantify uncertainty in the computed statistic by constructing an envelope (gray shaded region) via an appropriate bootstrap procedure. In the left panel, corresponding to a region selected inside the tumor, it is evident that the statistic deviates below the envelope at small scales indicating spatial correlation among the shapes of nuclei that manifests itself through less variation than expected under independence. At larger scales, the statistic deviates above the envelope indicating that the nuclei shapes are spatially correlated with more variation than expected if they were independent. In the right panel, which considers a region within non-tumor tissue, the nuclei shapes are uncorrelated in space since the statistic never deviates outside of the envelope. Thus, spatial correlation patterns for cell nuclei shapes within tumor and non-tumor tissues are distinct.

The correlation between coordinates of the point process and its marks may be studied using variants of Ripley's K function developed for homogeneous [Ripley, 1976] and

inhomogeneous Baddeley et al. [2000] point processes. In the presence of marks, weights associated with each pair of spatial locations are estimated using appropriate *test functions* [Illian et al., 2008], which include the mark correlation function [Stoyan and Stoyan, 1994] (summarizing average numerical similarity), the mark variogram [Cressie, 2015] (summarizing average variability as a function of spatial distance), and Moran’s I [Moran, 1950] (summarizing covariance); see Illian et al. [2008] for a more comprehensive discussion. Depending on the chosen test function, one obtains different second-order summary measures. A variety of marks has been studied in the statistics literature, including integer-valued (qualitative, multi-type or multivariate) [Lotwick and Silverman, 1982, van Lieshout and Baddeley, 1999], real-valued [Schlather, 2001, Schlather et al., 2004, van Lieshout, 2006], or mixed integer- and real-valued [Eckardt and Mateu, 2019].

Extensions to function-valued marks appear in Comas et al. [2008, 2011, 2013]. Most recently, Ghorbani et al. [2021] introduced a comprehensive framework that incorporates functional data with auxiliary variables as marks, proposing a mark-weighted n th-order reduced moment measure (based on test functions) to generalize existing summary statistics. In particular, they considered a mark-weighted inhomogeneous K function to quantify spatial correlatedness of the marks, while adjusting for the varying intensity of the spatial point process. Other extensions consider multivariate function-valued marks Eckardt and Moradi [2024] and compositional marks Eckardt et al. [2025].

The shape setting considered in this paper lies outside the capabilities of the above works with functional marks. Accordingly, we

- (a) define a novel population mark-weighted K function that accounts for the shapes of cell nuclei outlines as functional marks with symmetries, and a corresponding estimator that is easy to compute and is unbiased under some conditions;
- (b) use the estimator to construct global and local envelope tests for spatial dependence of the marks, which help answer questions (i) and (ii) above.

In the present setting, any test function used to define a mark-weighted K function must be invariant to the symmetries of translation, rescaling, rotation and reparameterization. The key issue lies in ensuring invariance to reparameterizations, which constitute an infinite-

dimensional Lie group. This is addressed by using an elastic Riemannian metric, invariant to reparameterizations, that then results in a nonlinear functional shape (or size-and-shape or orientation-and-shape) mark space. At the population level, the mark-weighted K function is defined using a product measure on the space of closed curves and desired symmetries, which at the sample level leads to a simple estimator.

The rest of this manuscript is organized as follows. Section 2 reviews elastic shape analysis of closed curves. Section 3 defines a marked spatial point process wherein marks are closed curves with desired symmetries. Section 4 introduces the shape, size-and-shape and orientation-and-shape mark-weighted K functions and provides a recipe for their estimation. Section 5 defines a formal hypothesis test of independence between marks and their spatial locations wherein the proposed mark-weighted K functions serve as test statistics. Section 6 presents simulations that assess the proposed framework. In Section 7, the framework is applied to cell nuclei outlines derived from histopathology images of breast cancer. Finally, we close with a brief discussion in Section 8. The supplement contains (i) a discussion of the reference measure on the mark space of shapes (Section 1), (ii) a derivation of the normalizing constant for the proposed mark-weighted K function (Section 2), (iii) a proof that the proposed K function estimator is unbiased (Section 3), (iv) a brief discussion of challenges associated with establishing consistency of the K function estimator (Section 4), (v) four additional histopathology images used in our analysis (Section 5), (vi) an additional result for cell nuclei from a region sampled in non-tumor tissue (Section 6), and (vii) code and data used to compute the results presented in Sections 6 and 7.

2 Comparing curves using elastic metrics

In this section, we briefly review the elastic framework for comparing planar parameterized curves, and refer to Srivastava and Klassen [2016] for details.

2.1 Curve representation and symmetries

Key to quantifying correlations between shapes of cell nuclei and their locations within the tissue is to systematically compare nuclei contours. This requires establishing cor-

respondence between points on two or more distinct curves. To this end, we represent cell nuclei as closed parameterized planar curves belonging to the set $\mathcal{B} := \{\beta : \mathbb{S}^1 \rightarrow \mathbb{R}^2 \mid \beta \text{ is absolutely continuous}\}$. The benefits of comparing objects given by planar outlines under a parameterized curve representation instead of just a subset of \mathbb{R}^2 are now well-understood [Chapters 6-7 Srivastava and Klassen, 2016].

As discussed earlier, the class of relevant symmetries associated with cell nuclei outlines will depend on many factors, including the TME and its location within the tissue. There are four transformations of a parameterized curve β that constitute symmetries pertaining to its shape: translation, scaling, rotation and reparameterization. Given a curve $\beta \in \mathcal{B}$, for every $t \in \mathbb{S}^1$, its translation is $\beta(t) + a$ where $a \in \mathbb{R}^2$, its scaling is $\sigma\beta(t)$ with $\sigma > 0$, its rotation is $O\beta(t)$, where $O \in SO(2) = \{O \in \mathbb{R}^{2 \times 2} \mid OO^\top = O^\top O = I_2, \det(O) = 1\}$; its reparameterization is the function composition $\beta \circ \gamma$, where $\gamma : \mathbb{S}^1 \rightarrow \mathbb{S}^1$ is a diffeomorphism of \mathbb{S}^1 , such that the image of $\beta \circ \gamma$ equals that of β for every such γ . Since reparameterizing β representing the outline of a cell nucleus preserves its image, it ought to be considered as a symmetry for curves, *irrespective* of any of the other factors that may influence the spatial configuration of the cell nuclei. In fact, the set $\Gamma := \{\gamma : \mathbb{S}^1 \rightarrow \mathbb{S}^1 \mid \gamma \text{ is a diffeomorphism}\}$ is an infinite-dimensional Lie group, with composition as the group operation. The shape symmetries then are represented by the product group, referred to as the *symmetry group*, $\mathbb{R}^2 \rtimes (\mathbb{R}_+ \times SO(2) \times \Gamma)$, with the semi-direct product action $((a, \sigma, O, \gamma), \beta) \mapsto \sigma O\beta \circ \gamma + a$. The order in which one scales, rotates and reparameterizes a curve is immaterial, since the individual symmetries commute. After these operations, a translation is applied.

The sizes of cell nuclei plausibly influence their spatial scatter, and it may be desirable then to not treat scale as a symmetry. In this case, the symmetry group is $\mathbb{R}^2 \rtimes (SO(2) \times \Gamma)$. Alternatively, the orientation of the nuclei might encode relevant information on infiltration of the tumor cell in a region of the tissue, and thus not nuisance; in this case, the symmetry group is $\mathbb{R}^2 \rtimes (\mathbb{R}_+ \times \Gamma)$. As such, we will consider three symmetry groups: (i) translation, scale, rotation and reparameterization, (ii) translation, rotation and reparameterization, and (iii) translation, scale and reparameterization. We generically denote a symmetry group by G , and refer to its elements as *shape symmetries*, since the image-preserving reparameterization group is always considered to be a symmetry. We denote the action of

a symmetry g on a curve β by $\beta * g$. Upon choosing G , the next step is to compare curves in a manner that is invariant to G . In other words, given two curves β_1 and β_2 , a distance d such that $d(\beta_1 * g, \beta_2 * g) = d(\beta_1, \beta_2)$ for every $g \in G$ is desired.

2.2 Elastic metric and the square-root velocity transform

The geometric framework [Joshi et al., 2007, Mio et al., 2007] based on the *elastic Riemannian metric* on \mathcal{B} provides a convenient set of tools to define a distance on curves compatible with the required symmetries. Definition of the metric requires a fair amount of preparatory notation, and thus, for ease of exposition, we refer the reader to Chapter 5 of Srivastava and Klassen [2016]. Summarily, the metric quantifies differences between two curves β_1 and β_2 by measuring the minimal deformation needed for bending, with weight $a > 0$, and stretching, with weight $b > 0$, whilst being invariant to reparameterization and rotation. This is a key requirement, since the parameterization of a curve, regardless of the chosen symmetry group, is nuisance. However, computing the distance using the elastic Riemannian metric for any choice of weights a, b is difficult. The situation is salvaged by the Square-Root Velocity (SRV) transform of a curve $\beta \in \mathcal{B}$, defined as

$$q(t) := \begin{cases} \frac{\dot{\beta}(t)}{\sqrt{|\dot{\beta}(t)|}} & \text{if } |\dot{\beta}(t)| \neq 0, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where $\dot{\beta}(t) = \frac{d\beta(t)}{dt}$ is the velocity vector. When $\beta \in \mathcal{B}$, its SRV transform q satisfies $\|q\|^2 := \int_{\mathbb{S}^1} |q(t)|^2 dt < \infty$, where $|\cdot|$ is the standard norm corresponding to the inner product $\langle \cdot, \cdot \rangle$ in \mathbb{R}^2 , and $\|\cdot\|$ is the norm on $\mathbb{L}^2(\mathbb{S}^1, \mathbb{R}^2)$ arising from the inner product $\langle \langle \cdot, \cdot \rangle \rangle$. Moreover, since β is a closed curve, a constraint imposed on q is $\int_{\mathbb{S}^1} \dot{\beta}(t) dt = \int_{\mathbb{S}^1} q(t)|q(t)| dt = 0$.

Let $\mathbb{L}^2(\mathbb{S}^1, \mathbb{R}^2)$ be the Hilbert space of square-integrable functions from the unit circle to the plane. The relevant function space is then $\mathcal{Q} := \{q \in \mathbb{L}^2(\mathbb{S}^1, \mathbb{R}^2) \mid \int_{\mathbb{S}^1} q(t)|q(t)| dt = 0\}$, which is a co-dimension two nonlinear submanifold of the infinite-dimensional Hilbert space $\mathbb{L}^2(\mathbb{S}^1, \mathbb{R}^2)$ [p.173 Srivastava and Klassen, 2016]. When using the weights $a = 1/4$ and $b = 1$ in the elastic metric on \mathcal{B} the elastic metric reduces to the standard \mathbb{L}^2 metric on \mathcal{Q} under the SRV transform so that the distance between q_1 and q_2 invariant to the chosen symmetry

group can be defined by using the usual \mathbb{L}^2 distance [Joshi et al., 2007].

2.3 Curve correspondence and distance under symmetries

Location/translation information of $\beta \in \mathcal{B}$ is lost when using its SRV transform q . When scale of q is a symmetry, then its normed version $q\|q\|^{-1}$ can be considered; rotation and reparameterization symmetries, however, need to be accounted for when comparing β_1 and β_2 via their SRV transforms. The action of a rotation O on β is identical to that on q . However, when β is reparametrized as $\beta \circ \gamma$ for $\gamma \in \Gamma$, the corresponding action of γ on its SRV transform q is given by $(q, \gamma) := (q \circ \gamma)\sqrt{\dot{\gamma}}$.

To establish correspondence between curves in \mathcal{B} invariant to a symmetry group G , we define a distance on the quotient \mathcal{Q}/G consisting of equivalence classes or orbits of q under the action of G . Curves invariant to their symmetries together are represented as points in quotient spaces, defined as follows.

- *Shape space* with $G = SO(2) \times \Gamma$:

$$\mathcal{Q}_{\text{sh}} := \{[q]_{\text{sh}}, q \in \mathcal{Q}\} \text{ with } [q]_{\text{sh}} = \{q * g := O(q\|q\|^{-1}, \gamma) \mid (O, \gamma) \in G\}.$$

- *Size-and-shape space* with $G = SO(2) \times \Gamma$:

$$\mathcal{Q}_{\text{sc-sh}} := \{[q]_{\text{sc-sh}}, q \in \mathcal{Q}\}, \text{ with } [q]_{\text{sc-sh}} = \{q * g := (Oq, \gamma) \mid (O, \gamma) \in G\}.$$

- *Orientation-and-shape space* with $G = \Gamma$:

$$\mathcal{Q}_{\text{ro-sh}} := \{[q]_{\text{ro-sh}}, q \in \mathcal{Q}\}, \text{ with } [q]_{\text{ro-sh}} = \{q * g := (q\|q\|^{-1}, \gamma) \mid \gamma \in G\}.$$

The shape of q is defined as its equivalence class under the action of $SO(2) \times \Gamma$, such that $q\|q\|^{-1}$ and $O(q\|q\|^{-1}, \gamma)$ have the same shape for any $(O, \gamma) \in SO(2) \times \Gamma$; size-and-shape and orientation-and-shape of q have analogous interpretations.

We now define distances on the three spaces, which will eventually enables us to quantify how cell nuclei and their symmetries influence their spatial configuration. Their definitions are based on the observation that, for every $(O, \gamma) \in SO(2) \times \Gamma$, $\|O(q, \gamma)\|^2 = \int_{\mathbb{S}^1} \left\langle Oq(\gamma(t))\sqrt{\dot{\gamma}(t)}, Oq(\gamma(t))\sqrt{\dot{\gamma}(t)} \right\rangle dt = \int_{\mathbb{S}^1} |q(u)|^2 du = \|q\|^2$, so that the rotation and parameterization groups act by isometries on \mathcal{Q} . This ensures that suitable distances may be defined on the shape, size-and-shape and orientation-and-shape spaces, as follows. The extrinsic shape distance between two curves $q_1, q_2 \in \mathcal{Q}$ is defined as $d_{\text{sh}}(q_1, q_2) =$

$\inf_{(O,\gamma) \in SO(2) \times \Gamma} \|q_1 \|q_1\|^{-1} - O(q_2 \|q_2\|^{-1}, \gamma)\|$, and is a distance on the shape space \mathcal{Q}_{sh} . In similar fashion, the size-and-shape distance is $d_{\text{sc-sh}}(q_1, q_2) = \inf_{(O,\gamma) \in SO(d) \times \Gamma} \|q_1 - O(q_2, \gamma)\|$, interpreted as a distance on the size-and-shape space $\mathcal{Q}_{\text{sc-sh}}$. Finally, the extrinsic orientation-and-shape distance is defined as $d_{\text{ro-sh}}(q_1, q_2) = \inf_{\gamma \in \Gamma} \|q_1 \|q_1\|^{-1} - (q_2 \|q_2\|^{-1}, \gamma)\|$, viewed as a distance on the orientation-and-shape space $\mathcal{Q}_{\text{ro-sh}}$. We refer to the three distances on curves generically as *shape distances*.

The distances d_{sh} and $d_{\text{ro-sh}}$ are referred to as extrinsic since they are based on the chordal distance on the infinite-dimensional sphere within \mathcal{Q} . Strictly speaking, d_{sh} , $d_{\text{sc-sh}}$, and $d_{\text{ro-sh}}$ are bonafide distance functions only if each equivalence class is closed; this is guaranteed by taking the closure of each equivalence by enlarging the group Γ to a monoid by including reparameterizations that are weakly increasing [Section 5.5 Srivastava and Klassen, 2016]. Computing the distances requires optimizing over $SO(2)$ and Γ . For a fixed γ , the solution to optimizing over $SO(2)$ is available in closed form via Procrustes analysis. For a fixed O , a dynamic programming algorithm, with an additional seed search over starting points on \mathbb{S}^1 , is used to optimize over Γ [Robinson, 2012]. Joint optimization over rotations and reparameterizations is carried out by alternating between the two.

A consequence of computing the distances above between a curve q_1 and another curve q_2 is the determination of the optimal element \hat{g} of the symmetry group G that best ‘matches’ q_1 to $q_2 * \hat{g}$. For example, if the shape distance $d_{\text{sh}}(q_1, q_2) = \|q_1 \|q_1\|^{-1} - \hat{O}(q_2 \|q_2\|^{-1}, \hat{\gamma})\|$ is realized at $(\hat{O}, \hat{\gamma})$, then, with respect to the shape distance, $q_1 \|q_1\|^{-1}$ is optimally aligned, or in shape correspondence, with $\hat{O}(q_2 \|q_2\|^{-1}, \hat{\gamma})$.

2.4 Computing means and geodesics

Higher-order moment measures for point processes with marks require both, pairwise and multiple comparison of marks. While the former is carried out for curves as marks by computing the shape distances above, a template is required to establish correspondence between multiple curves. A natural choice for a template is a Karcher mean.

Using the shape distances d_j , $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$, define the Karcher mean μ_j of a sample q_1, \dots, q_n of curves in \mathcal{Q} as the minimizer of $q \mapsto \operatorname{argmin}_{q \in \mathcal{Q}} \sum_{i=1}^n d_j(q, q_i)^2$. By construction, if μ_j is the Karcher mean with respect to d_j for $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$, then so

is $\mu_j * g$ for any g in the corresponding symmetry group. For example, $O(\mu_{\text{sh}}, \gamma)$ is a Karcher mean with respect to the shape distance d_{sh} for any $(O, \gamma) \in SO(2) \times \Gamma$ if μ_{sh} is. The Karcher mean is thus an equivalence class. In practice, computing the Karcher mean and multiple alignment of sample curves to the Karcher mean are performed simultaneously through an iterative algorithm. This algorithm alternates between aligning the sample curves to the current estimate of the Karcher mean and updating the mean based on the newly aligned curves. Ultimately, this process yields the Karcher mean and the sample shapes aligned to it. For further details, see Srivastava and Klassen [2016].

Upon determining the optimal symmetry \hat{g} when computing a distance d_j between two curves q_1 and q_2 , the geodesic between q_1 and q_2 is given by $(1 - \tau)q_1 + \tau(q_2 * \hat{g})$, $0 \leq \tau \leq 1$. For example, when using the size-and-shape distance $d_{\text{sc-sh}}$, suppose $d_{\text{sc-sh}}(q_1, O(q_2, \gamma))$ is minimized at $\hat{g} = (\hat{O}, \hat{\gamma}) \in SO(2) \times \Gamma$. The size-and-shape geodesic between q_1 and q_2 is given by $(1 - \tau)q_1 + \tau\hat{O}(q_2, \hat{\gamma})$, $0 \leq \tau \leq 1$; by construction, the midpoint $0.5q_1 + 0.5\hat{O}(q_2, \hat{\gamma})$ is then the Karcher mean $\mu_{\text{sc-sh}}$.

3 Point process with marks as curves with symmetries

The goal is to quantify the correlation between cell nuclei modulo their symmetries (e.g., shape, size-and-shape) and their locations within a histopathology image. To this end, we model the image as a realization of a marked point process in a spatial domain, with marks assuming values in the space of curves compatible with a chosen set of symmetries.

3.1 Spatial domain, mark space and reference measures

Consider $\mathcal{X} \subseteq \mathbb{R}^2$ equipped with the standard Euclidean norm $|\cdot|$ with Borel sets $\mathbb{B}(\mathcal{X})$. Given a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let $\Psi_{\text{gr}} = \{x_i\}_{i=1}^N$, $N \in \{0, 1, \dots, \infty\}$ on $(\mathcal{X}, \mathbb{B}(\mathcal{X}))$ be a point process, referred to as the *ground process*. Denote by G one of the symmetry groups: $SO(2) \times \Gamma$ for shape and size-and-shape, and Γ for orientation-and-shape.

We work with the space \mathcal{Q} containing SRV transformed curves from \mathcal{B} , and its quotients. Upon choosing a symmetry group G for the curves representing nuclei contours, there are two candidates for the mark space \mathcal{M} : (i) the product space $\mathcal{Q} \times G$, or (ii) a quotient

of \mathcal{Q} , such as the shape space \mathcal{Q}_{sh} . As we shall soon see, defining a second-order mark-weighted K function requires defining reference measures on \mathcal{X} and \mathcal{M} . When \mathcal{M} is a quotient of \mathcal{Q} containing equivalence classes of curves, a measure on \mathcal{Q} invariant to the group of symmetries is a natural choice for a reference measure. However, no such measure exists since each symmetry group under consideration is infinite-dimensional, on account of the Lie group Γ of reparameterizations of \mathbb{S}^1 . This presents significant difficulties in constructing a well-defined K function, and in its estimation. Our strategy is to thus consider $\mathcal{M} = \mathcal{Q}$, a Polish space (Section 1 in supplement), with Borel sets $\mathbb{B}(\mathcal{M})$, but introduce the necessary symmetries into the K function via a weight or test function.

At each point $x_i \in \mathcal{X}$ we attach a mark $(q_i * g_i) \in \mathcal{Q}$, resulting in the marked point process $\Psi = \{(x_i, (q_i * g_i))\}_{i=1}^N$ on $(\mathcal{X} \times \mathcal{M}, \mathbb{B}(\mathcal{X} \times \mathcal{M}))$. Let the spatial domain \mathcal{X} be equipped with the Lebesgue measure dx on \mathbb{R}^2 as its reference measure. The reference measure on \mathcal{M} is ν , defined as a pushforward of a product reference measure on $\mathcal{Q} \times G$ under the continuous group action map $(q, g) \mapsto q * g$; see Section 1 in supplement for a precise definition. This results in the product reference measure $\lambda := dx \otimes \nu$ on $\mathbb{B}(\mathcal{X} \times \mathcal{M})$. The symmetry g_i associated with curve q_i is interpretable only in a relative sense with respect to another curve – this is the central aspect of establishing correspondence. A consequence of this is seen in our definition of the K function, and its estimation, wherein a template is assumed.

3.2 Product density and pair correlation

Assume that the first-order product density ρ and the second-order product density $\rho^{(2)}$ of Ψ exist. Then, from the Campbell formula [Chiu et al., 2013], for a measurable function $f : \mathcal{X} \times \mathcal{M} \rightarrow [0, \infty)$,

$$\begin{aligned} & \mathbb{E} \left[\sum_{(x_1, q_1 * g_1), (x_2, q_2 * g_2) \in \Psi}^{\neq} f((x_1, q_1 * g_1), (x_2, q_2 * g_2)) \right] \\ &= \int_{(\mathcal{X} \times G)^2} f((x_1, q_1 * g_1), (x_2, q_2 * g_2)) \rho^{(2)}((x_1, q_1 * g_1), (x_2, q_2 * g_2)) \prod_{i=1}^2 d\lambda(x_i, q_i * g_i), \end{aligned} \tag{2}$$

where \neq represents summation over distinct tuples, and the measure $\prod_{i=1}^2 d\lambda(x_i, q_i * g_i) = \prod_{i=1}^2 dx(x_i) d\nu(q_i * g_i)$. The second-order product density $\rho^{(2)}$ may be expressed in terms of the second-order product density $\rho_{\text{gr}}^{(2)}$ of the ground process Ψ_{gr} as [Heinrich, 2012]

$$\rho^{(2)}((x_1, q_1 * g_1), (x_2, q_2 * g_2)) = J_{x_1, x_2}(q_1 * g_1, q_2 * g_2) \rho_{\text{gr}}^{(2)}(x_1, x_2), \quad (3)$$

where $J_{x_1, x_2}(q_1 * g_1, q_2 * g_2)$ represents the family of conditional joint densities of $(q_1 * g_1, q_2 * g_2)$ given pairs of locations (x_1, x_2) , where the densities are with respect to the reference measure ν . Indeed, families of conditional densities J_{x_1, \dots, x_k} exists for every $k \leq N$.

Denote by $\eta_{\text{gr}}(x_1, x_2) = \frac{\rho_{\text{gr}}^{(2)}(x_1, x_2)}{\rho_{\text{gr}}(x_1)\rho_{\text{gr}}(x_2)}$ the pair correlation function of the ground process Ψ_{gr} . Then, a more interpretable second-order measure of dependence within points in Ψ is its *pair correlation function*, defined as

$$\eta_{\Psi} = \frac{J_{x_1, x_2}}{J_{x_1} J_{x_2}} \eta_{\text{gr}}. \quad (4)$$

In (4), the first term represents the conditional joint density of the marks, providing insight into mark interactions given spatial locations. Meanwhile, η_{gr} reveals interactions between spatial points: values greater than 1 indicate clustering, values less than 1 suggest repulsion, and a value of 1 represents spatial independence. Although it is not a correlation in the traditional sense, this pair correlation function is convenient for examining second-order dependencies in spatial patterns.

We make the following three assumptions.

- A1.** The pair correlation function η_{Ψ} satisfies $\eta_{\Psi}((x_1, q_1 * g_1), (x_2, q_2 * g_2)) = \eta_{\Psi}((x_1 + v, q_1 * g_1), (x_2 + v, q_2 * g_2))$ for every $v \in \mathcal{X}$.
- A2.** The densities J_x are the same at all $x \in \mathcal{X}$, and the mark distribution is hence the same at all locations.
- A3.** The mark distribution equals the reference measure ν .

Under Assumption A1, the process Ψ is said to be second-order intensity-weighted stationary [Iftimi et al., 2019], which is a form of weak stationarity of Ψ . Further, dependence between points $(x_i, q_i * g_i)$ in Ψ depends only on the relative distances between their spatial locations

x_i , regardless of the marks. Assumption A1 is consistent with what is required to answer the motivational question (ii) in the Introduction. Assumption A2 offers a simplification when defining the K function. If, additionally, the marks are independent across spatial locations, the situation is referred to as *random labeling* [Ghorbani et al., 2021]. Assumption A3 simplifies computation of the K function, and will be discussed later.

4 Mark-weighted K function and estimation

We first define the population K function that represents a second-order dependence measure of the marked point process Ψ . The definition is based on the functional mark-weighted K function introduced by Ghorbani et al. [2021], suitably adapted to the setting of parameterized curves in the presence of a symmetry group G .

4.1 K function

The mark-weighted K function extends classical K functions by incorporating a test function that quantifies dependency in variations between marks at different locations. We consider three test functions $f : \mathcal{Q} \times \mathcal{Q} \rightarrow \mathbb{R}_{\geq 0}$ arising from the three distances d_{sh} , $d_{\text{sc-sh}}$, $d_{\text{ro-sh}}$ on curves, corresponding to three different symmetry groups, resulting in three K functions, which capture correlations between the spatial locations and nuclei shapes. The K function is obtained by weighting the second-order product densities at points within balls of radius $r > 0$ on the spatial domain \mathcal{X} by the test function.

Denote by $c_f := \mathbb{E}(f(q * g, q_1 * g_1) \mid q * g, q_1 * g_1 \perp X)$ the expectation of the test function f under random labeling, where the marks are independent of the locations. Let $B_x(r)$ be the open ball in \mathbb{R}^2 of radius $r > 0$ centered at a point x , and let \mathbb{I}_A denote the indicator function on the set A . Under Assumption A1 of stationarity of Ψ , a sensible mark-weighted K function may be defined as

$$K_f(r) = \frac{1}{|\mathcal{X}|c_f} \mathbb{E} \left[\sum_{(x, q * g) \in \Psi} \sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} \frac{f(q * g, q_1 * g_1)}{\rho(x, q * g)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \right],$$

where ρ is the first-order intensity of Ψ and $|\mathcal{X}|$ is the area of the spatial domain. The

test function \mathbf{f} summarizes dependency amongst the shape marks, and the K function thus represents the average value of \mathbf{f} for all pairs of marks separated by a distance r , normalized by the first-order product density to remove the effect of spatial intensity variations. The K function reflects how the similarity (or dissimilarity) of the marks depends on spatial proximity, allowing us to detect clustering, dispersion or spatial independence of the marks. Under random labeling, the normalizing constant is $c_{\mathbf{f}} = \mathbb{E}(\|q * g - m_G\|^2)$, where m_G is $\mathbb{E}(q * g)$ for $g \in G$ with respect to the measure ν on the mark space (Section 2 in the supplement). It is however difficult to estimate the product densities $\rho(x, q * g)$, especially since the mark space is infinite-dimensional and nonlinear.

We thus consider a simplified definition of the K function obtained by further making Assumptions A2 and A3: upon assuming a common mark distribution, from (3) we observe that if the mark distribution equals the reference measure, we get $\rho(x, q * g) = \rho_{\text{gr}}(x)$ for any $q \in \mathcal{Q}$ and $g \in G$, since the family of conditional densities $J_x(q * g) = J(q * g) \equiv 1$ for all $x \in \mathcal{X}$. Such an approach is also used in Iftimi et al. [2019], Ghorbani et al. [2021], D’Angelo et al. [2024]. Assumption A3 represents a practical solution given the infinite-dimensionality and non-Euclidean nature of the mark space \mathcal{M} . We thus consider a mark-weighted K function defined as follows.

Definition 1. Under the Assumptions A1-A3 of the process Ψ , the *mark-weighted K function* $r \mapsto K_{\mathbf{f}}(r)$, based on test function \mathbf{f} corresponding to a symmetry group G , is defined by

$$K_{\mathbf{f}}(r) = \frac{1}{|\mathcal{X}|c_{\mathbf{f}}} \mathbb{E} \left[\sum_{(x, q * g) \in \Psi} \sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} \frac{\mathbf{f}(q * g, q_1 * g_1)}{\rho_{\text{gr}}(x)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho_{\text{gr}}(x_1)} \right]. \quad (5)$$

The definition employs scaling by the intensity ρ_{gr} of the ground process to ensure that, under random labeling, $K_{\mathbf{f}}(r) = \pi r^2$, similar to what would be the case if the ground process Ψ_{gr} were to be a homogeneous Poisson point process on \mathcal{X} . This property of $K_{\mathbf{f}}$ under random labeling will be used in Section 5 when testing for the spatial dependence of the shape marks. We use three test functions \mathbf{f} based on the three curve distances:

$$\mathbf{f}(q * g, q_1 * g_1) := \frac{1}{2} d_j(q * \hat{g}, q_1 * \hat{g}_1)^2, \quad j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\},$$

where \hat{g} and \hat{g}_1 are the symmetries of q and q_1 that are in optimal correspondence with a template curve, say, μ_{temp} ; our definition of the mark-weighted K function implicitly then depends on the template μ_{temp} . The need for a template is due to the fact that the symmetries are interpretable only in a relative sense, and, moreover, is unavoidable in practice when working with quotient spaces under group actions. The factor of $\frac{1}{2}$ in \mathbf{f} cancels out a factor of 2 that arises in the derivation of the normalizing factor $c_{\mathbf{f}}$ (Section 2 in the supplement). Finally, the distance-based test functions may be interpreted as variograms that account for different symmetry groups.

4.2 Estimator for the K function and its computation

Given curves q_1, \dots, q_N , to estimate the K function in (5), we need to perform the following tasks: (i) choose a fixed template μ_{temp} or estimate one, (ii) estimate the unobserved shape symmetries g using the template ((i) and (ii) together help evaluate the test function \mathbf{f}), (iii) estimate the normalizing constant $c_{\mathbf{f}}$, (iv) estimate the intensity ρ_{gr} of the ground process Ψ_{gr} , and (v) correct for the boundary effect of $\mathcal{X} \subset \mathbb{R}^2$.

Two estimators of $K_{\mathbf{f}}$ may be considered. If the symmetries g , constant $c_{\mathbf{f}}$ and intensity ρ_{gr} are *known* or may be recovered exactly, then the estimator

$$\bar{K}_{\mathbf{f}}(r) = \frac{1}{|\mathcal{X}|c_{\mathbf{f}}} \sum_{(x, q * g) \in \Psi} \sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} \frac{w(x, x_1) \mathbf{f}(q * g, q_1 * g_1)}{\rho_{\text{gr}}(x) \rho_{\text{gr}}(x_1)} \mathbb{I}_{x \in \mathcal{X}} \mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}, \quad (6)$$

where w is an edge correction function satisfying $\int_{\mathcal{X}} \mathbb{I}\{(x + r) \in \mathcal{X}\} w(x, x + r) dx = |\mathcal{X}|$, is unbiased for $K_{\mathbf{f}}(r)$ for every $r > 0$; see Section 3 in the supplement for a proof. On the other hand if the symmetries are estimated consistently as sample size increases within \mathcal{X} (infill asymptotics), the estimator \hat{K} defined as $\bar{K}_{\mathbf{f}}$ with g , ground process intensity ρ_g and normalizing constant $c_{\mathbf{f}}$ replaced by estimators $\hat{g}, \hat{g}_1, \hat{\rho}_{\text{gr}}, \hat{\mathbf{f}}$ and $\hat{c}_{\mathbf{f}}$ in place of $g, g_1, \rho_{\text{gr}}, \mathbf{f}$ and $c_{\mathbf{f}}$, respectively, can be shown to be asymptotically unbiased with growing sample size within \mathcal{X} ; see Section 4 in the supplement for a discussion on consistency of $\hat{K}_{\mathbf{f}}$.

Examples of edge correction functions w satisfying the above requirement include the minus sampling edge correction $w_{\ominus}(x, x + y) = \frac{|\mathcal{X}| \mathbb{I}\{x \in \mathcal{X} \ominus B_0(|y|)\}}{|\mathcal{X} \ominus B_0(|y|)|}$, where \ominus denotes Minkowski

subtraction ($A \ominus B = \{x : x + A \subseteq B\}$), the translational edge correction $w_{\cap}(x, x + y) = \frac{|\mathcal{X}|}{|(\mathcal{X} + (x + y)) \cap (\mathcal{X} + x)|}$, and the isotropic edge correction $w(x, x + y) = \frac{2\pi|y|}{l(\partial B_x(|y|) \cap \mathcal{X})}$, where l is the length of the curve representing the boundary. Note that by an abuse of notation $|\cdot|$ denotes the area of a spatial domain as well as the Euclidean norm of a vector in \mathbb{R}^2 .

In practice, however, the symmetries g , constant c_f and intensity ρ_{gr} are unknown, and need to be estimated from the data. At a location x with mark $q * g$, for every x_1 with mark $q_1 * g_1$ within the ball of radius r around x , the symmetries g and g_1 are estimated by aligning, or establishing correspondence between, q and each q_1 with respect to the template. For the chosen test function f based on the distances d_j , $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$, we propose using the corresponding Karcher mean from Section 2.4 as an estimate of the template μ_{temp} . Thus, as described in Section 2.4, we compute the Karcher mean and perform joint alignment of q_1, \dots, q_N to the Karcher mean using an iterative procedure, resulting in the following estimates: (i) \hat{g} , \hat{g}_i for $i = 1, \dots, N$, (ii) Karcher mean μ_j , $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$ as estimate of the template μ_{temp} depending on the choice of the test function f , (iii) upon combining these, estimate $\hat{f}(q * \hat{g}, q_1 * \hat{g}_1)$ of the value of f .

Since the normalizing constant $c_f = \mathbb{E}(\|q * g - m_G\|^2)$, where m_G is $\mathbb{E}(q * g)$ for $g \in G$, we consider the estimate $\hat{c}_f = \frac{1}{N} \sum_{i=1}^N \|(q_i * \hat{g}_i) - \mu_j\|^2$, $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$, where μ_j is determined by the choice of f . Finally, we estimate the intensity ρ_{gr} of the ground process Ψ_{gr} using a Gaussian smoothing kernel. In later sections, we use the `density.ppp` function in the `spatstat` R package, where the bandwidth is selected following the approach of Cronie and Van Lieshout [2018] using the `bw.CvL` function in the same package.

5 Testing for spatial dependence of nuclei shapes

The mark-weighted K function and its estimate \hat{K} may be used to test if the nuclei shapes correlate with their spatial locations by comparing \hat{K} to the value of K under the null hypothesis of random labeling; under the null hypothesis, $K(r) = \pi r^2$ regardless of the intensity ρ_{gr} of the ground process Ψ_{gr} . This ensures that available tests for point processes, not necessarily with marks, may in principle be used.

The pipeline from a histopathology image to a realization Ψ of the marked point process

is noisy. It hence may be the case that spatial correlations between nuclei shapes manifest at different scales characterized by the radius r . We thus consider two types of tests: (i) a local envelope test [Ripley, 1977], which presents evidence for correlations for each radius r , and (ii) a global envelope test [Myllymäki et al., 2017], which provides evidence for correlations uniformly over all radii r . These tests were originally designed for point processes without marks. Adaptation to accommodate shape marks is straightforward owing to the aforementioned property of the definition of our K function.

The test statistic we use is based on a transformed version $\hat{L}(r) = \sqrt{\frac{\hat{K}(r)}{\pi}}$ of \hat{K} :

$$T(r) = \hat{L}(r) - L_{H_0}(r), \quad (7)$$

where $L_{H_0} = r$ if we wish to test for random labeling. The distribution of $T(r)$ is approximated via resampling under the null hypothesis of independence between marks and their locations, by permuting the marks while keeping their locations fixed. The test statistic with $L_{H_0} = r$ would result in a fairly conservative rejection region, since $L_{H_0} = r$ is sensible only when unknown components (e.g., symmetries, ground intensity ρ_{gr}) of the population K function in (5) are known *a priori*. Instead, as recommended by Grabarnik et al. [2011], we choose L_{H_0} as the \hat{L} obtained from the ground process when the test function $\mathbf{f}(\cdot, \cdot) \equiv 1$.

Once multiple samples of $T(r)$ have been obtained through resampling, for the local and global envelope tests, we provide, respectively:

- the proportion of radii for which the test statistic deviated from the 95% pointwise (in r) envelope generated using $s = 2499$ permutations [Ripley, 1977];
- the p-value derived from the global rank envelope test using the *extreme rank depth* measure on the samples with $s = 2499$ permutations [Myllymäki et al., 2017].

The global envelope test is implemented using the `global_envelope_test` function in the package `GET` in R. By carrying out the two types of tests, we note that if spatial dependence is asserted by the global test, further probing of dependence structures at local scales based on different values of r may be carried out by using the local test, along the lines of the method recently proposed by D’Angelo et al. [2024].

6 Simulations

We evaluate the effectiveness in detecting spatial dependence among shape, size-and-shape or orientation-and-shape marks based on the proposed mark-weighted K functions using simulations. In all cases, spatial dependence among marks is imposed via inter-point distances. We first simulate realizations of a homogeneous Poisson process (PP) over the domain $[0, 4] \times [0, 4]$ with an intensity of $\rho_{\text{gr}} = 8$, resulting in 128 locations on average. Computing for all experiments was performed on a node with 48 cores (Intel(R) Xeon(R) Platinum 8260 CPU @ 2.40GHz) with 187GB of memory. Estimation of a single shape mark-weighted K function in the simulations requires approximately 25 minutes on average (across ten replicates); no parallel computing was used. The computing time depends on (i) the number of curves in the spatial region of interest (average of 131), and (ii) discretization fineness for each curve (100).

6.1 Simulation of marks

We generate spatially correlated shapes at the sampled locations using the Fourier basis functions $\{\cos(lt), \sin(lt), l = 0, 1, 2, t \in [-\pi, \pi]\}$. Letting $\{b_j, j = 1, \dots, 6\}$ denote the basis functions, the radial distance as a function of the angle for curve i is defined as

$$\Gamma_i(t) = \sum_{j=1}^6 c_{ij} b_j(t), \quad (8)$$

where $\{c_{ij}, i = 1, \dots, N, j = 1, \dots, 6\}$ are Fourier coefficients. Each shape is then obtained by transforming the radial distance function to Cartesian coordinates: $(x_i(t) = \Gamma'_i(t) \cos(t), y_i(t) = \Gamma'_i(t) \sin(t))$, where $\Gamma'_i(t) = \Gamma_i(t) + |\min_t \Gamma_i(t)| + |\max_t \Gamma_i(t)|$ is an adjusted radial distance function that ensures that the resulting shapes are non-self-intersecting.

Spatial dependence in shape, rotation and scale is introduced independently. Shape dependence is incorporated through the Fourier coefficients c_{ij} , which are sampled from $\text{MVN}(\mathbf{0}_N^\top, C_{\text{mat}})$ to induce spatial dependence, and from $\text{MVN}(\mathbf{0}_N^\top, 0.1I_N + 0.9J_N)$ for the independent case. Here, MVN is short for the multivariate Normal distribution, $\mathbf{0}_N^\top$ is an N -dimensional vector of zeros, C_{mat} is the Matérn covariance matrix with scale parameter

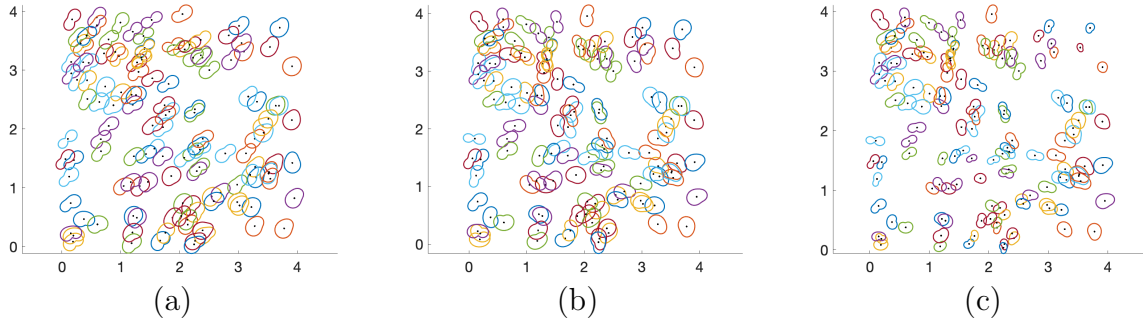


Figure 2: Simulated (a) spatially dependent shapes, (b) the same shapes as in (a) with spatially dependent orientations, and (c) the same shapes as in (b) with spatially dependent sizes.

of 0.5, smoothness parameter of 0.5 and range parameter of 2, I_N is an $N \times N$ identity matrix, and J_N is an $N \times N$ matrix of ones. To induce dependence in the shapes' orientations, we generate two unit vectors e_{i1} and e_{i2} independently from $\text{MVN}(0_N^\top, C_{\text{mat}})$, and compute the rotation angle ϕ_i between them; as before, C_{mat} is the Matérn covariance with the same parameter values. In the independent case, ϕ_i is sampled independently from $\text{Unif}(0, 2\pi)$. To model scale dependence, the scale factor σ_i is sampled from the correlated uniform distribution on $[0.2, 1.2]$, by transforming a sample from $\text{MVN}(1_N^\top, C_{\text{mat}})$. For the independent case, the Matérn covariance is replaced by the identity matrix I_N . The resulting shapes, with additional rotations and rescalings, are given by $\beta_i(t) = (x_i^*(t) = \sigma_i[x_i(t) \cos(\phi_i) - y_i(t) \sin(\phi_i)], y_i^*(t) = \sigma_i[x_i(t) \sin(\phi_i) + y_i(t) \cos(\phi_i)])$.

Figure 2 illustrates the sequence of steps taken to generate one realization of different types of marked point processes. Panel (a) shows a shape marked point process, where nearby curves tend to have similar shapes. Panel (b) presents the same curves as in (a), but with additional spatially dependent orientations. Panel (c) shows the same curves as in (b), but after additional spatially dependent rescalings.

6.2 Results

We evaluated eight scenarios, varying spatial dependence in shape, orientation and size. We computed the proposed mark-weighted K functions using the shape, orientation-and-shape and size-and-shape test functions and performed global and local tests for each scenario. The results, based on 50 independent replicates, are summarized in Table 1. Columns in

(a)			(b)			(c)		
(i)	(ii)	(iii)	(i)	(ii)	(iii)	(i)	(ii)	(iii)
0	0	0	0.52 (0.31)	0.48 (0.30)	0.46 (0.28)	0.04 (0.10)	0.05 (0.08)	0.05 (0.09)
1	0	0	0.00 (0.00)	0.15 (0.16)	0.12 (0.22)	0.98 (0.03)	0.26 (0.23)	0.50 (0.31)
0	1	0	0.54 (0.25)	0.02 (0.10)	0.47 (0.29)	0.03 (0.08)	0.93 (0.19)	0.05 (0.09)
0	0	1	0.52 (0.31)	0.48 (0.30)	0.00 (0.01)	0.04 (0.10)	0.05 (0.08)	0.96 (0.08)
1	1	0	0.00 (0.00)	0.00 (0.00)	0.13 (0.18)	0.99 (0.02)	0.96 (0.12)	0.43 (0.34)
1	0	1	0.00 (0.00)	0.15 (0.16)	0.00 (0.00)	0.98 (0.03)	0.26 (0.23)	0.99 (0.02)
0	1	1	0.54 (0.25)	0.02 (0.10)	0.00 (0.00)	0.03 (0.08)	0.93 (0.19)	0.95 (0.06)
1	1	1	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.99 (0.02)	0.96 (0.12)	1.00 (0.01)

Table 1: (a) Spatial dependence structure: ‘0’=spatial independence, ‘1’=spatial dependence; (i) shape, (ii) orientation, (iii) size. Results, across 50 replicates, of (b) global rank envelope tests with $\alpha = 0.05$ (average p-values), and (c) local envelope tests (average proportions of radii for which test statistics deviated from pointwise 95% envelope): (i) shape, (ii) orientation-and-shape, (iii) size-and-shape; standard deviations reported in parentheses.

(a) report the spatial dependence structure in the (i) shape, (ii) orientation and (iii) size of the marks, with ‘0’ corresponding to independence and ‘1’ to dependence. Columns in (b) report the average p-values from the global rank envelope tests (using extreme rank depth), at a significance level of $\alpha = 0.05$, based on the (i) shape, (ii) orientation-and-shape and (iii) size-and-shape mark-weighted K functions. Similarly, columns in (c) report average proportions of radii r for which the (i) shape, (ii) orientation-and-shape and (iii) size-and-shape mark-weighted K functions deviated from the pointwise 95% envelope. In both cases, standard deviations are presented in parentheses. We describe the results obtained from the global tests in detail and note that the pointwise tests resulted in very similar patterns.

In the independent mark scenario (shape=orientation=size=0), the global rank envelope tests based on all three mark-weighted K functions resulted in non-significant and fairly high p-values. On the other hand, when the marks’ shape, orientation and size were all spatially dependent (shape=orientation=size=1), all three global tests resulted in extremely small p-values; this result is reassuring as it represents the scenario with the strongest spatial dependence signal.

When only the marks’ shape exhibited spatial dependence (shape=1, orientation=size=0), the shape mark-weighted K function was able to detect this structure: the global test resulted in near zero p-values across all replicates. In contrast, global tests based on orientation-and-shape and size-and-shape resulted in non-significant p-values on average

(0.15 and 0.12, respectively). These tests jointly assess variation in both shape and either orientation or size, and, since the marks' orientations and sizes were spatially independent in this setting, their inclusion diluted the signal, reducing the power to detect shape dependence alone. When only size exhibited spatial dependence (shape=orientation=0, size=1), only the size-and-shape-based global test resulted in p-values near zero. While the size-and-shape test function assesses variation in both shape and scale, it seems to have primarily captured scale dependence; the shape- and orientation-and-shape-based tests resulted in non-significant and fairly large p-values. Finally, when only the marks' orientations were spatially dependent (shape=size=0, orientation=1), the shape- and size-and-shape-based global tests resulted in large p-values of 0.54 and 0.47 on average, respectively. On the other hand, the test based on the orientation-and-shape mark-weighted K function effectively captured the spatial dependence in the marks' orientations, with an average p-value of 0.02. However, the standard deviation in this case was relatively large: 0.10. The relative nature of orientation with respect to shape may have contributed to the slightly larger average p-value and much larger standard deviation as compared to cases of spatially dependent shape (0.00 (0.00)) or size (0.00 (0.01)). To verify this, we conducted an additional study. Initially, the Fourier coefficients c_{ij} in (8) were sampled from a multivariate normal distribution with the covariance matrix $(1 - \tau)I_N + \tau J_N$ with $\tau = 0.9$. As the coefficient on J_N decreases, the c_{ij} become more variable, increasing (independent) shape variation. Repeating our analysis using the orientation-and-shape mark-weighted K function with the coefficient reduced to 0.8, 0.7 and 0.6, the average p-values increased to 0.05 (0.17), 0.07 (0.19) and 0.08 (0.22), respectively. Similarly, the average proportion of radii r for which the mark-weighted K function deviated from the pointwise envelope decreased from 0.93 to 0.85, 0.80 and 0.74, respectively. These results confirm that accuracy of estimating relative orientations is influenced by the extent of shape heterogeneity in the sample.

When two out of three sources of variation, i.e., shape and orientation, shape and size or orientation and size, were spatially dependent, the results of global tests based on the different mark-weighted K functions were consistent. For example, when both the shape and size of marks were spatially dependent, the global tests based on shape or size-and-shape mark-weighted K functions resulted in near zero p-values, while the

one based on orientation-and-shape resulted in a non-significant p-value of 0.15. Overall, the presented simulation results demonstrate that the mark-weighted K function, when used with suitable test functions, effectively captures various forms of spatial dependence, allowing for assessment of random labeling in curves with symmetries.

7 Application to breast cancer data

7.1 Dataset and preliminary analysis

We examine two types of breast cancer known for poor prognosis: Human Epidermal Growth Factor Receptor 2-positive (HER2+) and Triple Negative (TNBC). We examined five tissue images with manually annotated tumor areas, reviewed by a panel of board certified breast pathologists, from the TIGER Challenge [2025] (Tumor Infiltrating Lymphocytes in Breast Cancer). For each tissue image, we randomly selected ten subregions within both tumor and non-tumor areas, using a fixed square window of $22500\mu m^2$. This resulted in a total of 50 subregions from tumor tissues and 50 subregions from non-tumor tissues. A representative annotated tissue image and the sampled subregions (yellow squares) selected for our analysis are shown in Figure 3. The tumor is outlined in black. The other four images are displayed in Section 5 in the supplement. The red and blue rectangles highlight the areas where the subregions were sampled. Nuclei outlines were extracted from each subregion using Qupath software [Bankhead et al., 2017]. Estimation of a single shape mark-weighted K function based on nuclei in a tumor region requires approximately 39 minutes on average (across ten regions; average number of curves per region: 185; discretization fineness: 100). Similarly, based on nuclei in a non-tumor region, estimation requires approximately 29 minutes on average (across ten regions; average number of curves: 134; discretization: 100); as in the simulations, no parallel computing was used.

We first examine the benefits of accounting for symmetries in cell nuclei outlines by computing geodesics between them as well as their Karcher means, as described in Section 2. In particular, we assess the symmetries in the cell nuclei as captured by the three distances d_j , $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$. Figure 4 compares two tumor cell nuclei via geodesic

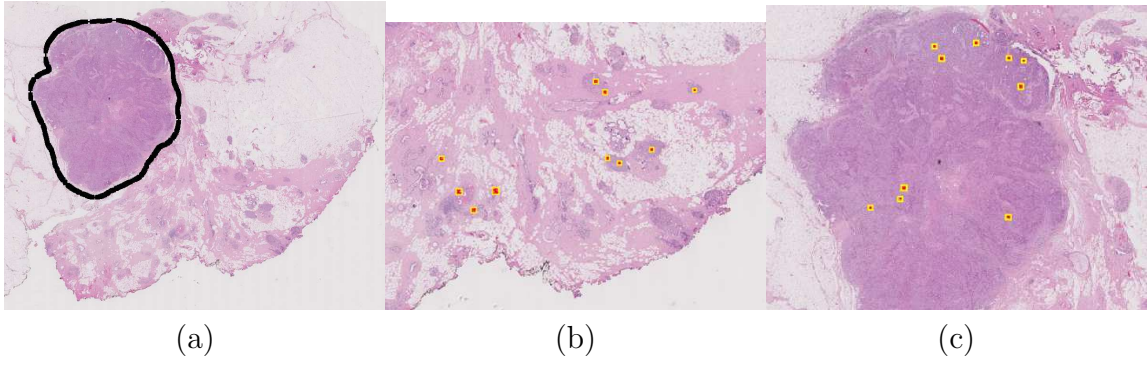


Figure 3: (a) Annotated tissue image: tumor area is marked by the black boundary. Ten regions (yellow squares) were randomly sampled within the (b) non-tumor area and (c) tumor area. The red (tumor) and blue (non-tumor) rectangles highlight the areas where regions were sampled.

paths and distances. The two nuclei differ in shape, size and relative orientation, as depicted in the left panel. Panels (a)-(c) illustrate comparisons from different perspectives. In (a), the geodesic path between the two nuclei outlines considers shape only (removing rotation and parameterization variations). Since rotation and parameterization are treated as symmetries, the shape distance is computed after the blue nucleus is optimally rotated and reparameterized with respect to the red one, and both shapes are scaled to have unit length. When orientation is included in the analysis, as shown in (b), the blue nucleus is optimally reparameterized with respect to the red one while maintaining its original orientation. In this case, both shapes are still scaled to unit length since size is also a symmetry. In contrast, panel (c) shows the size-and-shape geodesic wherein the blue nucleus was again optimally rotated and reparameterized with respect to the red one. However, the original sizes of the blue and red nuclei are preserved. When comparing the two nuclei, the geodesic distances vary depending on which factors are accounted for in the analysis. The shape distance between the two nuclei is 0.3079 (panel (a)). In (b), the distance, which now considers both orientation and shape, increases to 0.3763. This larger distance reflects the additional deformation required to account for the natural orientation of the shapes. Lastly, the size-and-shape distance in (c) is larger (2.5286) due to the difference in the natural sizes of the two shapes.

Figure 5(a) shows five tumor cell nuclei. Their resulting Karcher means as well as the optimally aligned shapes are shown in panels (b)-(d). In (d), the Karcher mean size-and-shape and all samples retain their natural sizes, as reflected by the axis values, while size

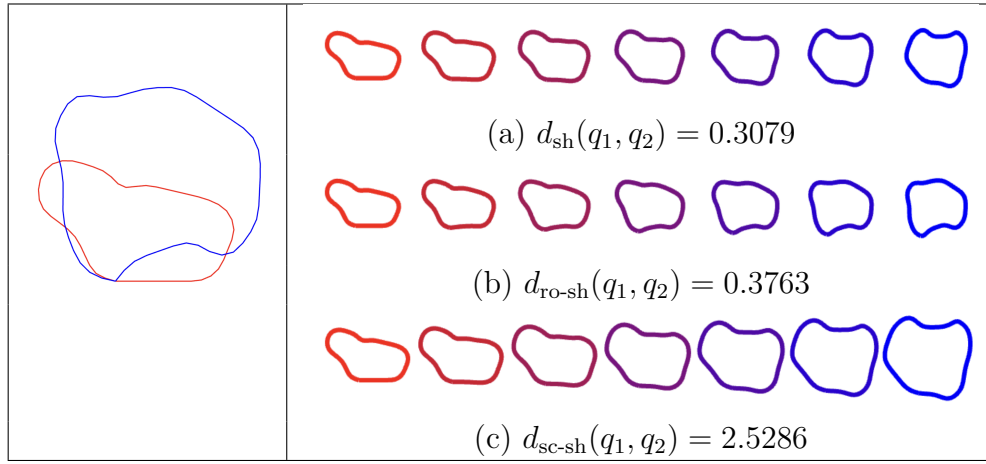


Figure 4: Left: Two nuclei boundaries from tumor area. Right: Geodesic paths on the (b) shape, (c) orientation-and-shape, and (d) size-and-shape spaces, with corresponding distances.

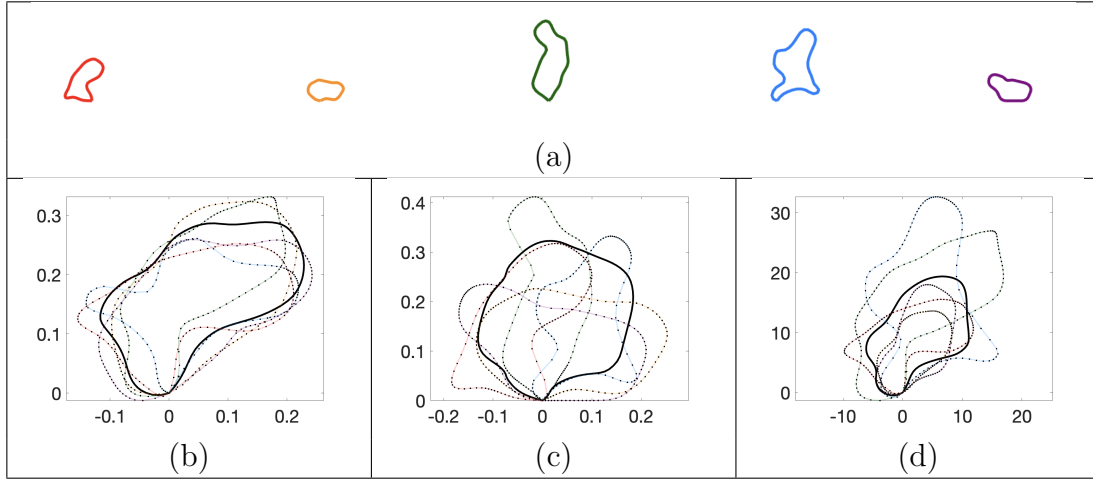


Figure 5: (a) Five nuclei outlines from tumor area. Karcher means (solid black) along with aligned sample nuclei (same colors as in (a)) based on (b) shape, (c) orientation-and-shape and (d) size-and-shape.

variation in (b) and (c) was removed. In both panels (b) and (d), the nuclei are aligned to their Karcher means by optimizing over both rotations and reparameterizations. In contrast, in panel (c), the nuclei retain their natural orientation, with alignment achieved solely by optimizing over reparameterizations. Overall, the results presented in Figures 4 and 5 show how the mark-weighted K function captures cell nuclei shape variation. First, the symmetries are accounted for through alignment of the nuclei to an appropriate Karcher mean. Then, variation is captured via appropriate test functions defined using the three distances d_j , $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$. Interpretation is aided by the ability to visualize geodesic paths between cell nuclei, whose lengths correspond to the computed distances.

7.2 Second-order analysis and practical implications

For each sampled region, we first estimated mark-weighted K functions, based on shape, orientation-and-shape or size-and-shape test functions. The mark-weighted K functions were then used to carry out (i) global rank envelope tests (using extreme rank depth) at a significance level of $\alpha = 0.05$, and (ii) local tests using 95% pointwise envelopes. P-values from the global tests are shown in the first row of Figure 6. The second row reports the proportions (converted to %) of radii r for which the K functions deviated from the pointwise envelope. In all cases, we also show the boxplots of the p-values or proportions. Panels (a)-(c) correspond to results for shape, orientation-and-shape and size-and-shape, respectively. Interesting patterns immediately emerge when examining these plots. In (a) and (c), it is evident that there is strong evidence of spatial dependence among the nuclei shapes and size-and-shapes for many more regions within tumor tissue as compared to non-tumor tissue. For normal tissue regions, the shape mark-weighted K function deviated from the pointwise envelope for small proportions radii: most between 0% and 20%. On the other hand, for tumor tissue regions, the proportions are much higher. This is further supported by the p-values from the global tests, which are much lower overall for tumor tissue regions (often lower than $\alpha = 0.05$). This pattern is even stronger when the size-and-shape mark-weighted K function is employed. This indicates that cell nuclei shapes and size-and-shapes are much more correlated with their spatial locations within tumor tissue than within non-tumor tissue. It further appears that cell nuclei within tumor and non-tumor tissues exhibit similar (strength of) spatial correlation in terms of orientation-and-shape.

Figure 7 shows how spatial dependence based on various types of variation emerges in a randomly chosen region within tumor tissue. Panel (a) presents the nuclei outlines. Panels (b)-(d) show the shape, orientation-and-shape and size-and-shape mark-weighted K functions, respectively, along with their 95% pointwise envelopes. Panels (e)-(g) display the corresponding test statistics used for the global tests along with their 95% global envelopes. In panel (e), it is evident that nuclei shapes do not exhibit spatial correlation up to approximately $r = 20$, after which they become more variable than expected under random labeling (deviation above the global envelope). In panel (g), the nuclei size-and-

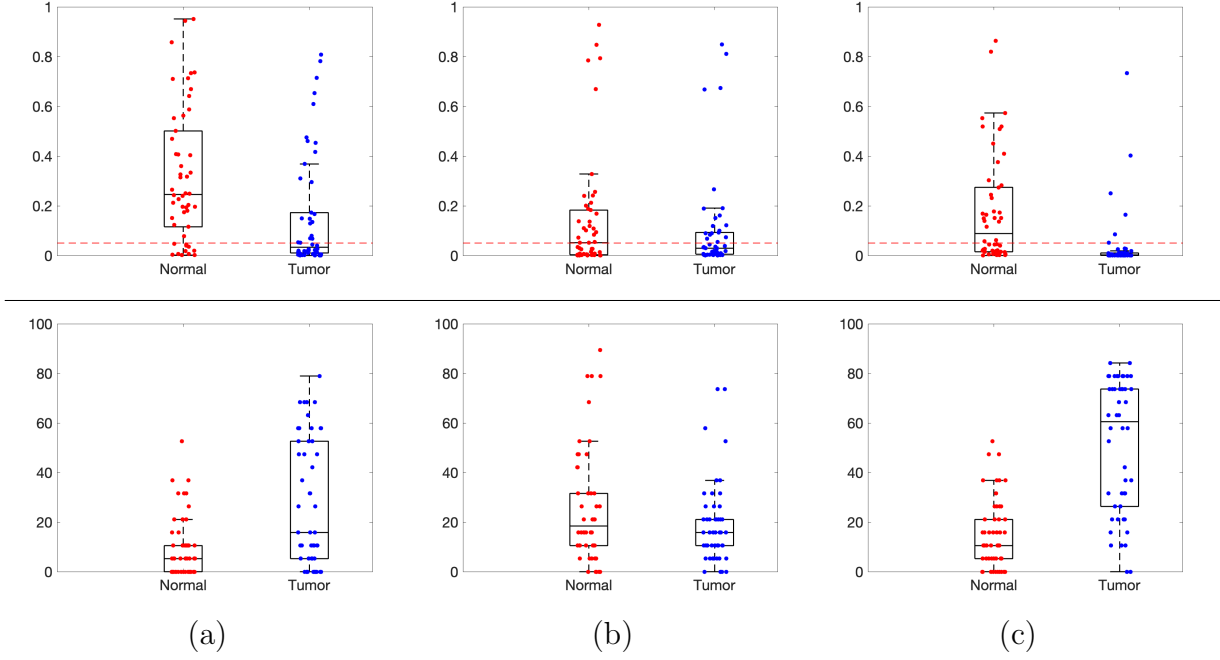


Figure 6: Boxplots of p-values from the global rank envelope test (top), and proportions of radii for which the mark-weighted K functions deviated outside the 95% pointwise envelope (bottom), based on the nuclei (a) shapes, (b) orientation-and-shapes, and (c) size-and-shapes. The red dashed line represents the significance level $\alpha = 0.05$ used in the global test. Points within each boxplot were jittered along the x -axis for clarity.

shapes exhibit smaller variability than expected by random labeling for radii close to $r = 10$. However, after approximately $r = 20$, they are more variable than expected under random labeling. In panel (f), we do not see evidence against random labeling based on the nuclei orientation-and-shapes. We observed very similar spatial dependence patterns in the other regions that were sampled within tumor tissues. For regions sampled in non-tumor tissues, the mark-weighted K functions (shape, orientation-and-shape or size-and-shape) did not consistently deviate from the pointwise or global envelopes; see Section 6 in the supplement.

The spatial dependency patterns for shape and size-and-shape aligned with expected tumor cell behavior. The mark-weighted K functions consistently remained within the expected range (envelope) for non-tumor cell nuclei. On the other hand, particularly based on size-and-shape, tumor cell nuclei showed significant deviations: values dropped below the envelope for radii of 6-13 μm , then rose above the envelope starting at a median radius of approximately 15.5 μm (the median radius was computed across all regions sampled within tumor tissues). Based on shape alone, the deviations above the envelope started

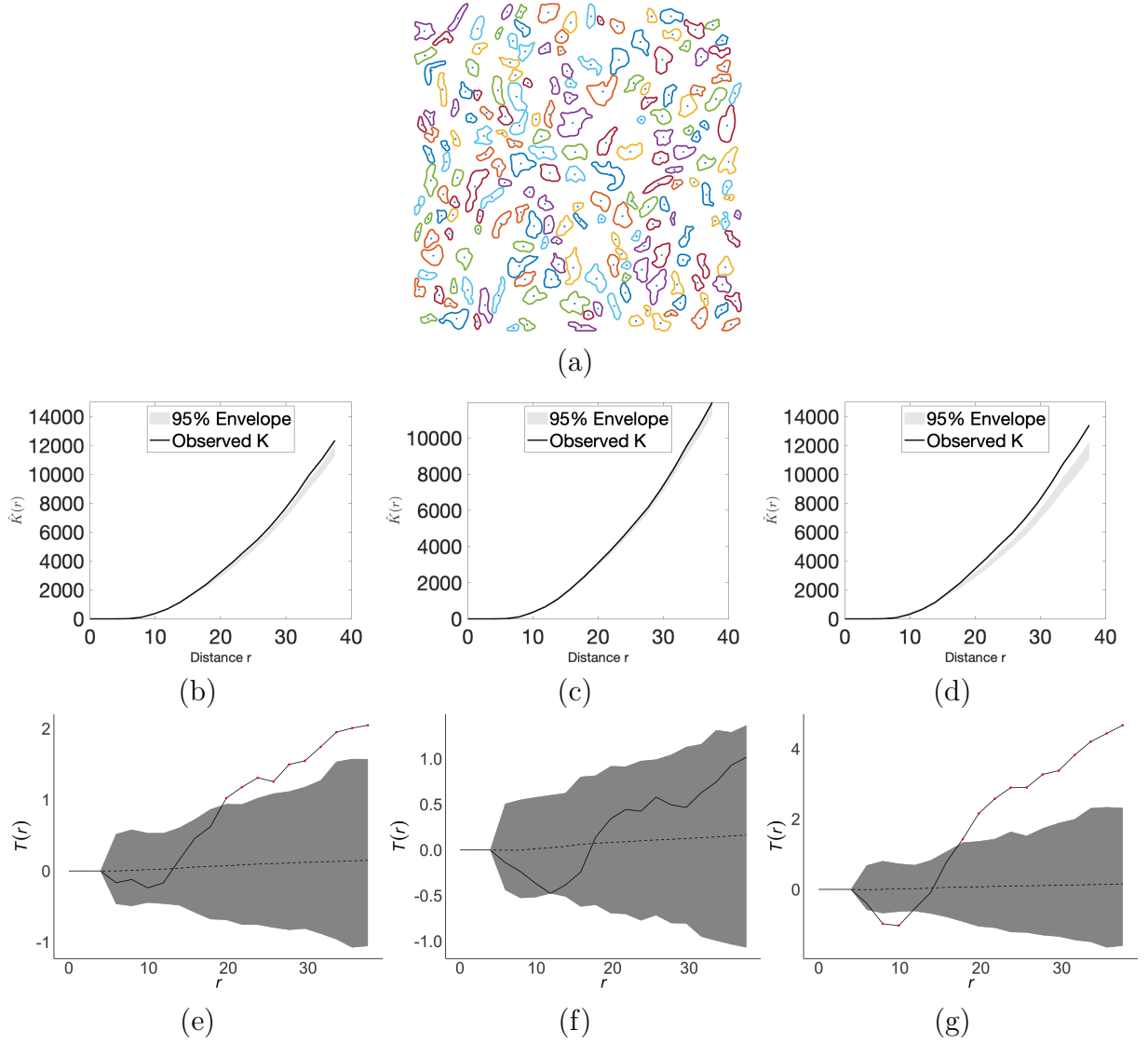


Figure 7: (a) A region within tumor tissue. Mark-weighted K and L functions along with their pointwise and global envelopes, respectively, estimated using (b)&(e) shape, (c)&(f) orientation-and-shape, and (d)&(g) size-and-shape test functions.

at a median radius of approximately $17.5 \mu m$. For reference, patient-derived circulating breast cancer cells have a median diameter of $12.4 \mu m$ [Mendelaar et al., 2021]; note that the cell nuclei are smaller in size. The deviation of the size-and-shape mark-weighted K functions below the envelope for small radii likely reflects tumor cells' distinct spatial organization. Normal breast cells arrange in tubular formations around ducts, whereas tumor cells form tight clusters of genetically similar 'clones' [Greaves and Maley, 2012]. This clustering of morphologically similar cells at short distances produces the observed reduction in K function values, as neighboring cells exhibit less variation than would be

expected by random labeling.

For larger radii, where clustering effects diminish, the significantly elevated K function values might reflect tumor heterogeneity – the diverse shapes and/or sizes of cancer cells within the same tumor. Unlike normal cells that maintain consistent nuclear morphology due to genetic stability, tumor cells accumulate different mutations and experience varying local conditions, creating a mixed population with markedly different appearances. This morphological diversity among neighboring cells results in more variation than expected under random labeling [Alizadeh et al., 2015, Zardavas et al., 2015, Marusyk et al., 2012].

8 Concluding remarks

Motivated by the problem of uncovering and testing for spatial dependence patterns among cell nuclei shapes in histopathology images of cancer, we propose a marked point process framework based on a mark-weighted K function statistic to capture shape correlations among cell nuclei. The framework uses the elastic shape analysis of curves and ensures invariance to desired symmetries in the underlying data. Simulations confirm the effectiveness of the novel K functions in capturing relevant spatial correlations for planar closed curves, and application of the proposed framework to histopathology images of breast cancer reveals distinct spatial dependence patterns in their shape and size-and-shape across tumor and non-tumor regions. Our findings corroborate existing hypotheses in the literature on breast cancer tumor growth and heterogeneity, as well as cell proliferation.

We have identified several directions for future work. First, we will build on the framework of Ghorbani et al. [2021] to quantify spatial relationships among shape marks modulated by auxiliary variables in a latent space. Second, we will extend the framework to apply to shapes of outlines of three-dimensional objects, i.e., surfaces. Cell nuclei are inherently three-dimensional and this will allow us to capture full information about their shape. Third, we will try to establish consistency of the proposed mark-weighted K functions when the reparameterization group is restricted to a finite-dimensional family; see discussion in Section 4 in the supplement. Finally, we plan to apply the proposed methods to other types of cancer, e.g., prostate, to try to understand what properties of cell nuclei

affect their spatial distribution. This has the potential to reveal different types of tumor heterogeneity in different contexts.

Acknowledgments. This work was supported in part by NIH R37-CA214955 (SK and KB), NSF DMS-2015374 and EPSRC EP/Z003377/1 (KB), and NSF DMS-2413747 (SK).

References

- AA Alizadeh, V Aranda, A Bardelli, C Blanpain, C Bock, C Borowski, C Caldas, A Califano, M Doherty, M Elsner, M Esteller, R Fitzgerald, JO Korb, P Lichter, CE Mason, N Navin, D Pe'er, K Polyak, CW Roberts, L Siu, A Snyder, H Stower, C Swanton, RG Verhaak, JC Zenklusen, J Zuber, and J Zucman-Rossi. Toward understanding and exploiting tumor heterogeneity. *Nature Medicine*, 21(8):846–853, 2015.
- AJ Baddeley, J Møller, and R Waagepetersen. Non- and semi-parametric estimation of interaction in inhomogeneous point patterns. *Statistica Neerlandica*, 54(3):329–350, 2000.
- P Bankhead, MB Loughrey, JA Fernández, Y Dombrowski, DG McArt, PD Dunne, S McQuaid, RT Gray, LJ Murray, HG Coleman, JA James, M Salto-Tellez, and Hamilton PW. Qupath: Open source software for digital pathology image analysis. *Scientific Reports*, 7(1):1–7, 2017.
- S Barua, L Solis, ER Parra, N Uraoka, M Jiang, H Wang, J Rodriguez-Canales, I Wistuba, A Maitra, S Sen, and A Rao. A functional spatial analysis platform for discovery of immunological interactions predictive of low-grade to high-grade transition of pancreatic intraductal papillary mucinous neoplasms. *Cancer Informatics*, 17:1176935118782880, 2018.
- AH Beck, AR Sangoi, S Leung, RJ Marinelli, TO Nielsen, MJ Van De Vijver, RB West, M Van De Rijn, and D Koller. Systematic analysis of breast cancer morphology uncovers stromal features associated with survival. *Science Translational Medicine*, 3(108), 2011.
- HJG. Bloom and WW Richardson. Histological grading and prognosis in breast cancer;

- a study of 1409 cases of which 359 have been followed for 15 years. *British Journal of Cancer*, 11:359–377, 1957.
- SN Chiu, D Stoyan, WS Kendall, and J Mecke. *Stochastic Geometry and its Applications*. John Wiley & Sons, 2013.
- C Comas, P Delicado, and J Mateu. Analysing spatial point patterns with associated functional data. In *Statistics for Spatio-temporal Modelling*, pages 157–163, 2008.
- C Comas, P Delicado, and J Mateu. A second order approach to analyse spatial point patterns with functional marks. *Test*, 20(3):503–523, 2011.
- C Comas, L Mehtätalo, and J Miina. Analysing space-time tree interdependencies based on individual tree growth functions. *Stochastic Environmental Research and Risk Assessment*, 27:1673–1681, 2013.
- N Cressie. *Statistics for Spatial Data*. John Wiley & Sons, 2015.
- O Cronie and MNM Van Lieshout. A non-model-based approach to bandwidth selection for kernel estimators of spatial intensity functions. *Biometrika*, 105(2):455–462, 2018.
- CE de Andrea, AS Petrilli, R Jesus-Garcia, LF Bleggi-Torres, and MTS Alves. Large and round tumor nuclei in osteosarcoma: Good clinical outcome. *International Journal of Clinical and Experimental Pathology*, 4(2):169, 2011.
- N D’Angelo, G Adelfio, J Mateu, and O Cronie. Local inhomogeneous weighted summary statistics for marked point processes. *Journal of Computational and Graphical Statistics*, 33(2):588–602, 2024.
- M Eckardt and J Mateu. Analysing multivariate spatial point processes with continuous marks: A graphical modelling approach. *International Statistical Review*, 87(1):44–67, 2019.
- M Eckardt and M Moradi. Marked spatial point processes: Current state and extensions to point processes on linear networks. *Journal of Agricultural, Biological and Environmental Statistics*, 29(2):346–378, 2024.

- Matthias Eckardt, Mari Myllymäki, and Sonja Greven. On spatial point processes with composition-valued marks. *arXiv:2501.16049*, 2025.
- CW Elston and IO Ellis. Pathological prognostic factors in breast cancer. I. The value of histological grade in breast cancer: Experience from a large study with long-term follow-up. *Histopathology*, 19(5):403–410, 1991.
- EG Fischer. Nuclear morphology and the biology of cancer cells. *Acta Cytologica*, 64(6): 511–519, 2020.
- M Ghorbani, O Cronie, J Mateu, and J Yu. Functional marked point processes: A natural structure to unify spatio-temporal frameworks and to analyse dependent functional data. *Test*, 30(3):529–568, 2021.
- P Grabarnik, M Myllymäki, and D Stoyan. Correct testing of mark independence for marked point patterns. *Ecological Modelling*, 222(23-24):3888–3894, 2011.
- M Greaves and CC Maley. Clonal evolution in cancer. *Nature*, 481(7381):306–313, 2012.
- D Hanahan and LM Coussens. Accessories to the crime: Functions of cells recruited to the tumor microenvironment. *Cancer Cell*, 21(3):309–322, 2012.
- D Hanahan and RA Weinberg. Hallmarks of cancer: The next generation. *Cell*, 144(5): 646–674, 2011.
- L Heinrich. Asymptotic methods in statistics of random point processes. In *Stochastic Geometry, Spatial Statistics and Random Fields: Asymptotic Methods*, pages 115–150. Springer, 2012.
- A Iftimi, O Cronie, and F Montes. Second-order analysis of marked inhomogeneous spatiotemporal point processes: Applications to earthquake data. *Scandinavian Journal of Statistics*, 46(3):661–685, 2019.
- J Illian, A Penttinen, H Stoyan, and D Stoyan. *Statistical Analysis and Modelling of Spatial Point Patterns*. John Wiley & Sons, 2008.

- SH Joshi, E Klassen, A Srivastava, and IH Jermyn. A novel representation for Riemannian analysis of elastic curves in \mathbb{R}^n . In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–7, 2007.
- JN Kather, J Krisam, P Charoentong, T Luedde, E Herpel, C-A Weis, T Gaiser, A Marx, NA Valous, D Ferber, L Jansen, CC Reyes-Aldasoro, I Zörnig, D Jäger, H Brenner, J Chang-Claude, M Hoffmeister, and N Halama. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS Medicine*, 16(1):e1002730, 2019.
- HW Lotwick and BW Silverman. Methods for analysing spatial processes of several types of points. *Journal of the Royal Statistical Society, Series*, 44(3):406–413, 1982.
- C Lu, K Bera, X Wang, P Prasanna, J Xu, A Janowczyk, N Beig, M Yang, P Fu, J Lewis, H Choi, RA Schmid, S Berezowska, K Schalper, D Rimm, V Velcheti, and A Madhabushi. A prognostic model for overall survival of patients with early-stage non-small cell lung cancer: A multicentre, retrospective study. *The Lancet*, 2(11):e594–e606, 2020.
- A Marusyk, V Almendro, and K Polyak. Intra-tumour heterogeneity: A looking glass for cancer? *Nature Reviews Cancer*, 12(5):323–334, 2012.
- P Mendelaar, J Kraan, M Van, LL Zeune, LWMM Terstappen, E Oomen-de Hoop, JWM Martens, and S Sleijfer. Defining the dimensions of circulating tumor cells in a large series of breast, prostate, colon, and bladder cancer patients. *Molecular Oncology*, 15(1):116–125, 2021.
- W Mio, A Srivastava, and SH Joshi. On shape of plane elastic curves. *International Journal of Computer Vision*, 73:307–324, 2007.
- PAP Moran. Notes on continuous stochastic phenomena. *Biometrika*, 37(1/2):17–23, 1950.
- GF Murphy, AW Partin, SJ Maygarden, and JL Mohler. Nuclear shape analysis for assessment of prognosis in renal cell carcinoma. *Journal of Urology*, 143(6):1103–1107, 1990.

- M Myllymäki, T Mrkvička, P Grabarnik, H Seijo, and U Hahn. Global envelope tests for spatial processes. *Journal of the Royal Statistical Society, Series B*, 79(2):381–404, 2017.
- R Nafe, K Franz, W Schlote, and B Schneider. Morphology of tumor cell nuclei is significantly related with survival time of patients with glioblastomas. *Clinical Cancer Research*, 11(6):2141–2148, 2005a.
- R Nafe, W Schlote, and B Schneider. Histomorphometry of tumour cell nuclei in astrocytomas using shape analysis, densitometry and topometric analysis. *Neuropathology and Applied Neurobiology*, 31(1):34–44, 2005b.
- BD Ripley. The second-order analysis of stationary point processes. *Journal of Applied Probability*, 13(2):255–266, 1976.
- BD Ripley. Modelling spatial patterns. *Journal of the Royal Statistical Society, Series B*, 39(2):172–192, 1977.
- DT Robinson. *Functional Data Analysis and Partial Shape Matching in the Square Root Velocity Framework*. The Florida State University, 2012.
- R Sali, Y Jiang, A Attaranzadeh, B Holmes, and R Li. Morphological diversity of cancer cells predicts prognosis across tumor types. *Journal of the National Cancer Institute*, 116(4):555–564, 2024.
- M Schlather. On the second-order characteristics of marked point processes. *Bernoulli*, 7(1):99–117, 2001.
- M Schlather, PJ Ribeiro Jr, and PJ Diggle. Detecting dependence between marks and locations of marked point processes. *Journal of the Royal Statistical Society, Series B*, 66(1):79–93, 2004.
- I Singh and TP Lele. Nuclear morphological abnormalities in cancer: A search for unifying mechanisms. In *Nuclear, Chromosomal, and Genomic Architecture in Biology and Medicine*, pages 443–467. Springer, 2022.

- S Singh, S Raman, E Caserta, G Leone, M Ostrowski, J Rittscher, and R Machiraju. Analysis of spatial variation of nuclear morphology in tissue microenvironments. In *IEEE International Symposium on Biomedical Imaging*, pages 1293–1296, 2010.
- A Srivastava and EP Klassen. *Functional and Shape Data Analysis*. Springer, 2016.
- D Stoyan and H Stoyan. *Fractals, Random Shapes and Point Fields: Methods of Geometrical Statistics*. John Wiley & Sons, 1994.
- TIGER Challenge. Tiger challenge: Tumor imaging genetics resource, 2025. URL <https://tiger.grand-challenge.org/>.
- T Tsujikawa, J Mitsuda, H Ogi, A Miyagawa-Hayashino, E Konishi, K Itoh, and S Hirano. Prognostic significance of spatial immune profiles in human solid cancers. *Cancer Science*, 111(10):3426–3434, 2020.
- C van Doijeweert, PJ van Diest, and IO Ellis. Grading of invasive breast carcinoma: The way forward. *Virchows Archiv*, 480(1):33–43, 2022.
- MNM van Lieshout. A J-function for marked point patterns. *Annals of the Institute of Statistical Mathematics*, 58(2):235–259, 2006.
- MNM van Lieshout and AJ Baddeley. Indices of dependence between types in multivariate point patterns. *Scandinavian Journal of Statistics*, 26(4):511–532, 1999.
- X Wang, C Barrera, K Bera, VS Viswanathan, S Azarianpour-Esfahani, C Koyuncu, P Velu, MD Feldman, M Yang, P Fu, KA Schalper, H Mahdi, C Lu, V Velcheti, and A Madhabushi. Spatial interplay patterns of cancer nuclei and tumor-infiltrating lymphocytes (TILs) predict clinical benefit for immune checkpoint inhibitors. *Science Advances*, 8(22):eabn3966, 2022.
- D Zardavas, A Irrthum, C Swanton, and M Piccart. Clinical management of breast cancer heterogeneity. *Nature Review: Clinical Oncology*, 12(7):381–394, 2015.

Supplement for second-order spatial analysis of shapes of tumor cell nuclei

Ye Jin Choi, Sebastian Kurtek

Department of Statistics, The Ohio State University

Simeng Zhu

James Cancer Center, The Ohio State University

Karthik Bharath

School of Mathematical Sciences, University of Nottingham

November 13, 2025

1 Reference measure on the mark space of shapes

The space \mathcal{Q} is a codimension two submanifold of $\mathbb{L}^2(\mathbb{S}^1, \mathbb{R}^2)$, and a Polish space. Let ζ_1 be the restriction of a non-degenerate Gaussian measure on $\mathbb{L}^2(\mathbb{S}^1, \mathbb{R}^2)$ to \mathcal{Q} . Consider first only the reparametrization group Γ as the symmetry. The reparametrization group Γ may be prescribed an appropriate topology under which it is a Polish space [Cohen, 2017]. Let ζ_2 be a quasi-invariant σ -finite measure on Γ [see, e.g., Section 11.5 Bogachev, 2010]. The group action $\Phi : \mathcal{Q} \times \Gamma \rightarrow \mathcal{Q}$, $(q, \gamma) \rightarrow (q \circ \gamma)\sqrt{\gamma}$ is continuous. Thus, the pushforward $\nu := \zeta_1 \times \zeta_2 \circ \Phi^{-1}$ of the product measure, also known as the convolution measure, on $\mathcal{Q} \times \Gamma$ is well-defined, with support on a closed subset of \mathcal{Q} .

The rotation group $SO(2)$ is a compact Lie group with a single connected component, and Polish. Let ζ_3 be the uniform Haar measure on $SO(2)$. The action $(Oq, \gamma) \rightarrow O(q \circ \gamma)\sqrt{\gamma}$ is also continuous, which ensures that the pushforward of the product measure $\zeta_1 \times \zeta_2 \times \zeta_3$ under the group action results similarly in a well-defined measure with support on a closed subset of \mathcal{Q} .

2 Normalizing constant in population K function

We show that, under random labeling, the normalizing constant takes the form $c_f = \mathbb{E}(\|q * g - m_G\|^2)$, where $m_G = \mathbb{E}(q * g)$ and $g \in G$.

$$\begin{aligned}
c_f &:= \mathbb{E}[\mathbf{f}(q * g, q_1 * g_1) \mid (q * g, q_1 * g_1) \perp X] \\
&= \mathbb{E}\left(\frac{1}{2} \|q * g - q_1 * g_1\|^2\right) \\
&= \frac{1}{2} \mathbb{E}(\|(q * g - m_G) + (m_G - q_1 * g_1)\|^2) \\
&= \mathbb{E}(\|q * g - m_G\|^2) + \int_{\mathbb{S}^1} \mathbb{E}[(q * g)(t) - m_G(t)] [(q_1 * g_1)(t) - m_G(t)] dt \\
&= \mathbb{E}(\|q * g - m_G\|^2)
\end{aligned}$$

The fourth equality follows from Fubini's theorem with dt as the arc length measure. The final equality holds because the marks $q * g$ and $q_1 * g_1$ are independent under the assumption of random labeling.

3 Unbiasedness of the estimator \bar{K}_f

The estimator \bar{K}_f in Equation (6) in the main manuscript is unbiased when the symmetries g are *known* or are exactly recovered. To see this, consider first the following intermediate result, which is a special case of Lemma 1 in Ghorbani et al. [2021] with an additional normalizing factor c_f . Let $\kappa_f(u_1, u_2)$ denote the 2nd order intensity-reweighted t -correlation measure (Definition 8 in Ghorbani et al. [2021]), defined by

$$\kappa_f(u_1, u_2) = \int_{\mathcal{M}^2} \frac{J_{u_1, u_2}}{J_{u_1} J_{u_2}} \mathbf{f}(q_1 * g_1, q_2 * g_2) \prod_{i=1}^2 d\nu(q_i * g_i).$$

Then, the K function can be expressed as

$$K_f(r) = \frac{1}{c_f} \int_{B_0(r)} \kappa_f(0, u_1) K_{\text{gr}}(du_1), \tag{1}$$

where K_{gr} is the 2nd order reduced moment measure of the ground process.

Let $\mathbb{E}^{l(x, (q, \gamma))}$ denote the expectation under the reduced Palm measure $P^{l(x, (q, \gamma))}$, which is the conditional distribution of a functional marked point process given that it has a point at $(x, (q * g))$. Then, we can express $c_f K_f(r)$ as

$$\begin{aligned} c_f K_f(r) &= \frac{1}{|\mathcal{X}|} \mathbb{E} \left[\sum_{(x, (q * g)) \in \Psi} \sum_{(x_1, (q_1 * g_1)) \in \Psi \setminus \{(x, (q * g))\}} \frac{f(q * g, q_1 * g_1)}{\rho(x, q * g)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \right] \\ &= \frac{1}{|\mathcal{X}|} \int_{\mathcal{X} \times \mathcal{M}} \mathbb{E}^{l(x, q * g)} \left[\sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} \frac{f(q * g, q_1 * g_1)}{\rho(x, q * g)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \right] d\lambda(x, q * g). \end{aligned}$$

The last equality is due to the Campbell-Mecke formula. Further, $\mathbb{E}^{l(x, q * g)}(\cdot)$ can be expressed as

$$\begin{aligned} &\mathbb{E}^{l(x, q * g)} \left[\sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} \frac{f(q * g, q_1 * g_1)}{\rho(x, q * g)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \right] \\ &= \int_{\mathcal{X} \times \mathcal{M}} \frac{f(q * g, q_1 * g_1)}{\rho(x, q * g)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \rho^{(2)}((x, q * g), (x_1, q_1 * g_1)) d\lambda(x_1, q_1 * g_1) \\ &= \int_{\mathcal{X} \times \mathcal{M}} f(q * g, q_1 * g_1) \mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\} \eta_{\Psi}((x, q * g), (x_1, q_1 * g_1)) d\lambda(x_1, q_1 * g_1) \\ &\stackrel{\Lambda 1}{=} \int_{\mathcal{X} \times \mathcal{M}} f(q * g, q_1 * g_1) \mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\} \eta_{\Psi}((0, q * g), (x_1 - x, q_1 * g_1)) d\lambda(x_1, q_1 * g_1) \\ &\stackrel{u_1 = x_1 - x}{=} \int_{B_0(r) \times \mathcal{M}} f(q * g, q_1 * g_1) \eta_{\Psi}((0, q * g), (u_1, q_1 * g_1)) d\lambda(u_1, q_1 * g_1) \\ &= \int_{B_0(r) \times \mathcal{M}} f(q * g, q_1 * g_1) \frac{J_{0, u_1}(q * g, q_1 * g_1)}{J_0(q * g) J_{u_1}(q_1 * g_1)} \eta_{\text{gr}}(0, u_1) d\lambda(u_1, q_1 * g_1) \\ &= \int_{B_0(r)} \left(\int_{\mathcal{M}} f(q * g, q_1 * g_1) \frac{J_{0, u_1}(q * g, q_1 * g_1)}{J_0(q * g) J_{u_1}(q_1 * g_1)} d\nu(q * g) \right) \eta_{\text{gr}}(0, u_1) du(u_1). \end{aligned}$$

Thus, we have

$$\begin{aligned}
& c_f K_f(r) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{X} \times \mathcal{M}} \mathbb{E}^{!(x, q * g)} \left[\sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} \frac{f(q * g, q_1 * g_1)}{\rho(x, q * g)} \frac{\mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \right] d\lambda(x, q * g) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{X} \times \mathcal{M}} \int_{B_0(r)} \left(\int_{\mathcal{M}} f(q * g, q_1 * g_1) \frac{J_{0, u_1}(q * g, q_1 * g_1)}{J_0(q * g) J_{u_1}(q_1 * g_1)} d\nu(q * g) \right) \times \\
&\quad \times \eta_{\text{gr}}(0, u_1) du(u_1) d\lambda(u, q * g) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{X}} \int_{B_0(r)} \left(\int_{\mathcal{M}^2} f(q * g, q_1 * g_1) \frac{J_{0, u_1}(q * g, q_1 * g_1)}{J_0(q * g) J_{u_1}(q_1 * g_1)} d\nu(q * g) d\nu(q_1 * g_1) \right) \times \\
&\quad \times \eta_{\text{gr}}(0, u_1) du(u_1) du(u) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{X}} \int_{B_0(r)} \kappa_f(0, u_1) \eta_{\text{gr}}(0, u_1) du(u_1) du(u) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{X}} du(u) \int_{B_0(r)} \kappa_f(0, u_1) \eta_{\text{gr}}(0, u_1) du(u_1) \\
&= \int_{B_0(r)} \kappa_f(0, u_1) K_{\text{gr}}(du_1).
\end{aligned}$$

Then,

$$\begin{aligned}
& \mathbb{E} [c_f \bar{K}_f(r)] \\
&= \frac{1}{|\mathcal{X}|} \mathbb{E} \left[\sum_{(x, q * g) \in \Psi} \sum_{(x_1, q_1 * g_1) \in \Psi \setminus \{(x, q * g)\}} w(x, x_1) f(q * g, q_1 * g_1) \times \right. \\
&\quad \left. \times \frac{\mathbb{I}\{x \in \mathcal{X}\}}{\rho(x, q * g)} \frac{1\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \right] \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{M} \times \mathcal{X}} \int_{\mathcal{M} \times \mathcal{X}} w(x, x_1) f(q * g, q_1 * g_1) \frac{\mathbb{I}\{x \in \mathcal{X}\}}{\rho(x, q * g)} \frac{1\{x_1 \in \mathcal{X} \cap B_x(r)\}}{\rho(x_1, q_1 * g_1)} \times \\
&\quad \times \rho^{(2)}((x, q * g), (x_1, q_1 * g_1)) d\lambda(x, q * g) d\lambda(x_1, q_1 * g_1)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{M}} \int_{\mathcal{M}} \mathbf{f}(q * g, q_1 * g_1) \times \\
&\quad \times \left(\int_{\mathcal{X}} \int_{\mathcal{X} \cap B_x(r)} w(x, x_1) \frac{\rho^{(2)}((x, q * g), (x_1, q_1 * g_1))}{\rho(x, q * g) \rho(x_1, q_1 * g_1)} dx(x_1) dx(x) \right) \times \\
&\quad \times d\nu(q * g) d\nu(q_1 * g_1) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{M}} \int_{\mathcal{M}} \mathbf{f}(q * g, q_1 * g_1) \times \\
&\quad \times \left(\int_{\mathcal{X}} \int_{\mathcal{X} \cap B_x(r)} w(x, x_1) \eta_{\Psi}((x, q * g), (x_1, q_1 * g_1)) dx(x_1) dx(x) \right) \times \\
&\quad \times d\nu(q * g) d\nu(q_1 * g_1),
\end{aligned}$$

and by assumption A1 and Fubini's theorem, the inner expression satisfies

$$\begin{aligned}
&\int_{\mathcal{X}} \int_{\mathcal{X} \cap B_x(r)} w(x, x_1) \eta_{\Psi}((x, q * g), (x_1, q_1 * g_1)) dx(x_1) dx(x) \\
&\stackrel{\text{A1}}{=} \int_{\mathcal{X}} \int_{\mathcal{X}} \mathbb{I}\{x_1 \in \mathcal{X} \cap B_x(r)\} w(x, x_1) \eta_{\Psi}((0, q * g), (x_1 - x, q_1 * g_1)) dx(x_1) dx(x) \\
&\stackrel{u_1 = x_1 - x}{=} \int_{\mathcal{X}} \int_{\mathcal{X}} \mathbb{I}\{u_1 + x \in \mathcal{X} \cap B_x(r)\} w(x, x + u_1) dx(x) \times \\
&\quad \times \eta_{\Psi}((0, q * g), (u_1, q_1 * g_1)) du(u_1) \\
&= \int_{\mathcal{X}} \mathbb{I}\{u_1 \in B_0(r)\} \left(\int_{\mathcal{X}} \mathbb{I}\{u_1 + x \in \mathcal{X}\} w(x, x + u_1) dx(x) \right) \times \\
&\quad \times \eta_{\Psi}((0, q * g), (u_1, q_1 * g_1)) du(u_1) \\
&= |\mathcal{X}| \int_{B_0(r)} \eta_{\Psi}((0, q * g), (u_1, q_1 * g_1)) du(u_1),
\end{aligned}$$

because

$$\int_{\mathcal{X}} \mathbb{I}\{u_1 + x \in \mathcal{X}\} w(x, x + u_1) dx(x) = |\mathcal{X}|.$$

Hence, by Fubini's theorem and the intermediate result in (1), we have the following:

$$\begin{aligned}
& \frac{1}{|\mathcal{X}|} \int_{\mathcal{M}} \int_{\mathcal{M}} \mathbf{f}(q * g, q_1 * g_1) \left(\int_{\mathcal{X}} \int_{\mathcal{X} \cap B_x(r)} w(x, x_1) \eta_{\Psi}((x, q * g), (x_1, q_1 * g_1)) dx(x_1) dx(x) \right) \\
& \quad \times d\nu(q * g) d\nu(q_1 * g_1) \\
&= \frac{1}{|\mathcal{X}|} \int_{\mathcal{M}} \int_{\mathcal{M}} \mathbf{f}(q * g, q_1 * g_1) \left(|\mathcal{X}| \int_{B_0(r)} \eta_{\Psi}((0, q * g), (u_1, q_1 * g_1)) du(u_1) \right) \\
& \quad \times d\nu(q * g) d\nu(q_1 * g_1) \\
&= \int_{\mathcal{M}} \int_{\mathcal{M}} \mathbf{f}(q * g, q_1 * g_1) \left(\int_{B_0(r)} \frac{J_{0,u_1}(q * g, q_1 * g_1)}{J_0(q * g) J_{u_1}(q_1 * g_1)} \eta_{\text{gr}}(0, u_1) du(u_1) \right) \\
& \quad \times d\nu(q * g) d\nu(q_1 * g_1) \\
&\stackrel{\text{Fubini}}{=} \int_{B_0(r)} \left[\int_{\mathcal{M}} \int_{\mathcal{M}} \frac{J_{0,u_1}(q * g, q_1 * g_1)}{J_0(q * g) J_{u_1}(q_1 * g_1)} \mathbf{f}(q * g, q_1 * g_1) d\nu(q * g) d\nu(q_1 * g_1) \right] \\
& \quad \times \eta_{\text{gr}}(0, u_1) du(u_1) \\
&= \int_{B_0(r)} \kappa_{\mathbf{f}}(0, u_1) \eta_{\text{gr}}(0, u_1) du(u_1) \\
&= \int_{B_0(r)} \kappa_{\mathbf{f}}(0, u_1) K_{\text{gr}}(du_1) \\
&\stackrel{(1)}{=} c_{\mathbf{f}} K_{\mathbf{f}}(B_0(r)) \\
&\equiv c_{\mathbf{f}} K_{\mathbf{f}}(r)
\end{aligned}$$

4 Challenges in establishing consistency of \hat{K}

Pointwise consistency of the estimator \hat{K} depends on, for each location $x \in \mathcal{X}$, the exact recovery of the symmetries g, g_1 , and values of the ground intensity ρ_{gr} at x and x_1 ; then, continuity of the distance functions d_j , $j \in \{\text{sh}, \text{sc-sh}, \text{ro-sh}\}$ ensures convergence of $\hat{\mathbf{f}}$ to \mathbf{f} , and from the law of large numbers in Hilbert spaces (since an extrinsic metric $\|\cdot\|$ is used on \mathcal{Q}), $\hat{c}_{\mathbf{f}}$ also converges to $c_{\mathbf{f}}$ under random labeling.

However, exact recovery of symmetries is not possible in practice. The key issue lies with the infinite-dimensionality of the symmetry group due to the presence of the reparametrization group Γ ; this issue can be seen in the simpler case of functions in the presence of phase variation, wherein exact recovery is possible only when the dimensionality of phase varia-

tion is rendered finite-dimensional [Chakraborty and Panaretos, 2021, Kurtek et al., 2011]. It may thus be possible to prove consistency by assuming a finite-dimensional parametric family of reparameterizations, but we leave that for future work.

Ghorbani et al. [2021] considered marked point processes with marks in an infinite-dimensional linear vector space in the absence of any symmetries and invariances, and proved consistency of an estimator of a K function without the normalizing factor c_f . They assumed ergodicity of Ψ and a known ground intensity function. However, their proof techniques cannot be used in our setting involving symmetries.

5 Additional annotated tissue images

Figure 1 displays the other four histopathology images that were used for the analysis in Section 7 in the main article.

6 Result for a region within non-tumor tissue

Figure 2 presents results akin to Figure 6 in the main manuscript, but for a region sampled in non-tumor tissue. Here, we see that there is no strong evidence against random labeling based on shape, orientation-and-shape or size-and-shape. This observation was consistent for most regions sampled within non-tumor tissues, and supports the clinical interpretation of our results presented in Section 7 in the main article.

References

- VI Bogachev. *Differentiable measures and the Malliavin calculus*. Number 164. American Mathematical Society, 2010.
- A Chakraborty and VM Panaretos. Functional registration and local variations: Identifiability, rank, and tuning. *Bernoulli*, 27:1103–1130, 2021.
- MP Cohen. Polishability of some groups of interval and circle diffeomorphisms. *arXiv preprint arXiv:1709.04523*, 2017.

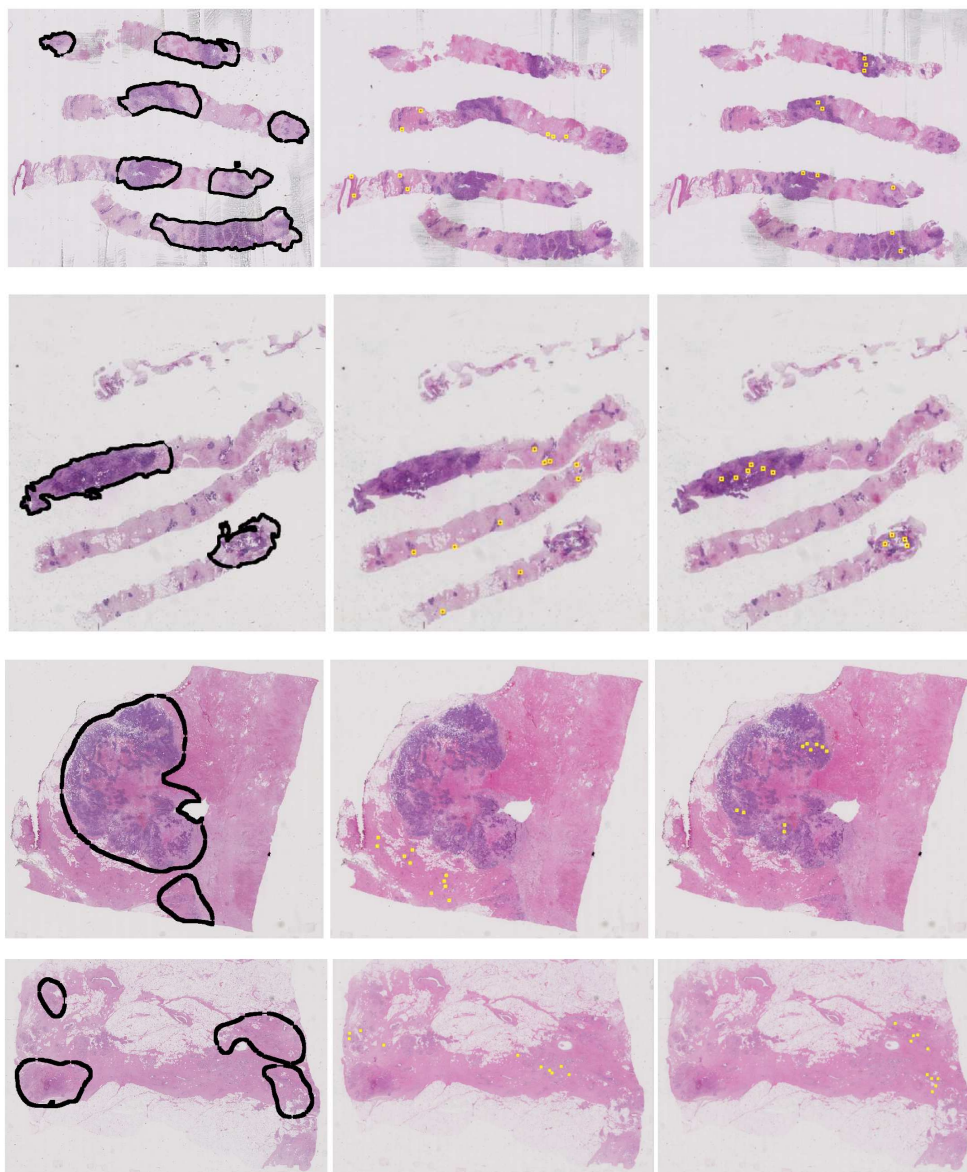


Figure 1: Left: Annotated tissue image: tumor area is marked by the black boundary. Ten regions (yellow squares) were randomly sampled within the non-tumor area (middle) and tumor area (right).

M Ghorbani, O Cronie, J Mateu, and J Yu. Functional marked point processes: A natural structure to unify spatio-temporal frameworks and to analyse dependent functional data. *Test*, 30(3):529–568, 2021.

S Kurtek, A Srivastava, and W Wu. Signal estimation under random time-warpings and nonlinear signal alignment. In *Advances in Neural Information Processing Systems*, volume 24, 2011.

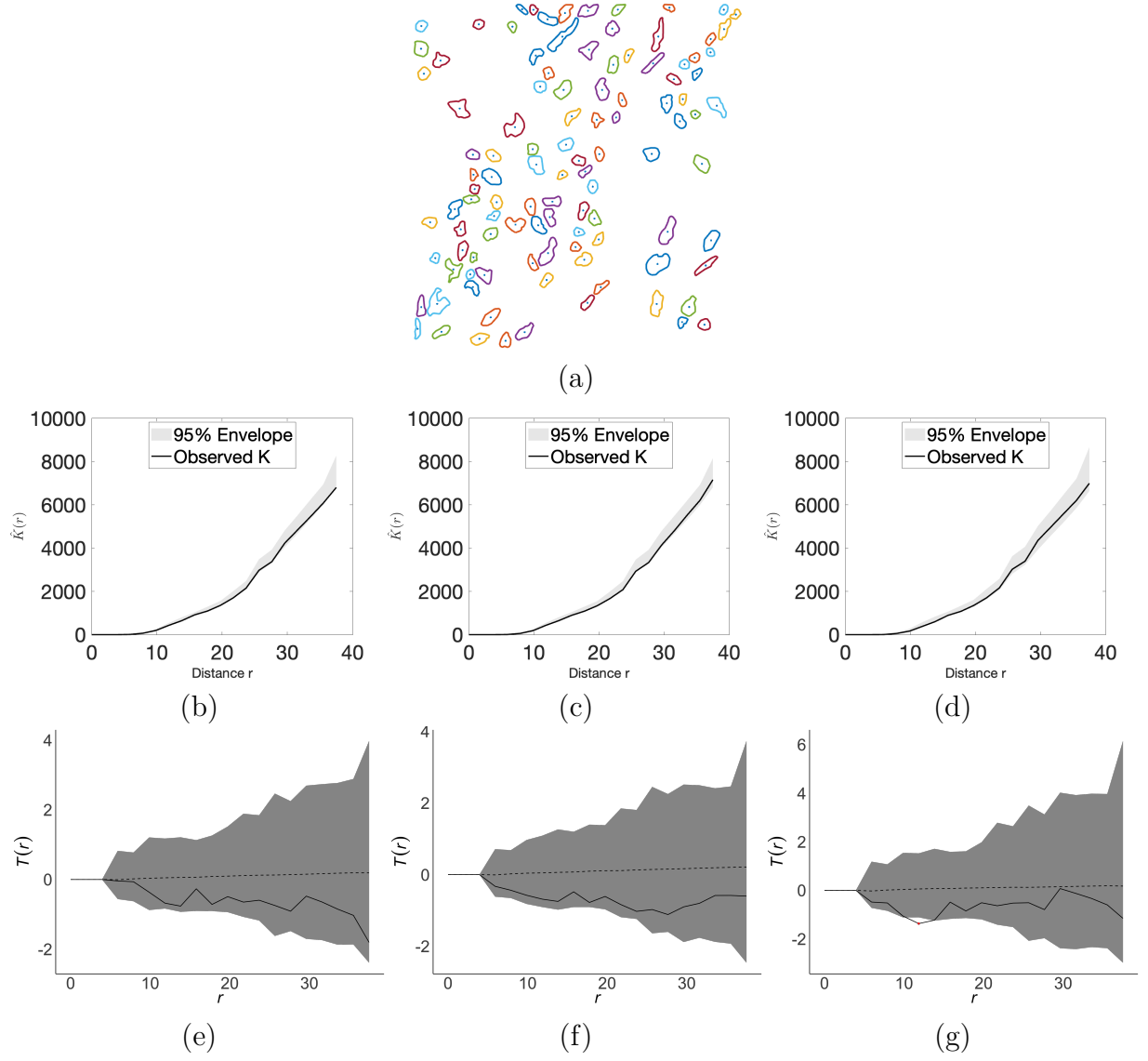


Figure 2: (a) A region within non-tumor tissue. Mark-weighted K and L functions along with their pointwise and global envelopes, respectively, estimated using (b)&(e) shape, (c)&(f) orientation-and-shape, and (d)&(g) size-and-shape test functions.