

Multivariate Functional Shape Data Analysis

1 INTRODUCTION

2 BACKGROUND AND METHODS

2.1 Background

Suppose that we observe an image dataset for n unrelated subjects. Without loss of generality, we focus on a compact set, denoted as $\mathcal{D} \subset \mathbb{R}^t$, which is general enough to cover curves ($t = 1$), contours ($t = 2$), and surfaces ($t = 3$). It is assumed that $\{\mathbf{d}_1, \dots, \mathbf{d}_{N_V}\}$ are N_V grid points (or vertices) on \mathcal{D} from the template file. Specifically, for the i -th subject, we observe a $J \times 1$ vector of shape measurements corresponding to each grid point \mathbf{d} , denoted as $\mathbf{y}_i(\mathbf{d}) = (y_{i1}(\mathbf{d}), \dots, y_{iJ}(\mathbf{d}))^T$, and a $p \times 1$ vector of covariates (e.g., age, gender, group information, and biological markers), denoted as \mathbf{x}_i with its first component being 1.

2.2 Multivariate Functional Shape Data Analysis (MFSDA)

The MFSDA is defined as

$$y_{ij}(\mathbf{d}) = \mathbf{x}_i^T \boldsymbol{\beta}_j(\mathbf{d}) + \eta_{ij}(\mathbf{d}) + \epsilon_{ij}(\mathbf{d}), \quad (1)$$

where $\boldsymbol{\beta}_j(\mathbf{d})$ is a $p \times 1$ vector of fixed effects, $\boldsymbol{\eta}_i(\mathbf{d}) = (\eta_{i1}(\mathbf{d}), \dots, \eta_{iJ}(\mathbf{d}))^T$ characterizes both subject-specific and location-specific variability, and $\boldsymbol{\epsilon}_i(\mathbf{d}) = (\epsilon_{i1}(\mathbf{d}), \dots, \epsilon_{iJ}(\mathbf{d}))^T$ are measurement errors. It is also assumed that $\boldsymbol{\eta}_i(\mathbf{d})$ and $\boldsymbol{\epsilon}_i(\mathbf{d})$ are mutually independent and identical

copies of $\text{SP}(\mathbf{0}, \Sigma_\eta)$ and $\text{SP}(\mathbf{0}, \Sigma_\epsilon)$, respectively, where $\text{SP}(\boldsymbol{\mu}, \Sigma)$ denotes a stochastic process vector with mean function $\boldsymbol{\mu}(\mathbf{d})$ and covariance function $\Sigma(\mathbf{d}, \mathbf{d}')$. Moreover, $\Sigma_\epsilon(\mathbf{d}, \mathbf{d}')$ takes the form of $\Omega_\epsilon(\mathbf{d})\mathbf{1}(\mathbf{d} = \mathbf{d}')$, where $\Omega_\epsilon(\mathbf{d})$ is a nonnegative function of \mathbf{d} and $\mathbf{1}(\cdot)$ is an indicator function of an event. Compared with the standard linear regression model, MFSDA explicitly accounts for spatial smoothness, spatial correlation, and the low-dimensional representation of functional shape responses [Zhu et al., 2014, 2011].

2.3 Hypothesis Testing

Under model (1), we start with a hypothesis testing problem on $\beta_j(\mathbf{d})$, $j = 1, \dots, J$, to investigate whether there is statistically significant morphological difference caused by some covariate of interest or linear combination of covariates of interest:

$$H_0 : \mathbf{C}\beta(\mathbf{d}) = \mathbf{0} \quad \text{v.s.} \quad H_1 : \mathbf{C}\beta(\mathbf{d}) \neq \mathbf{0} \text{ for each } \mathbf{d}, \quad (2)$$

where $\mathbf{C} = \mathbf{I}_J \otimes \mathbf{a}^T$, $\beta(\mathbf{d}) = \text{vec}([\beta_1(\mathbf{d}), \dots, \beta_J(\mathbf{d})])$, and \mathbf{a} is a $p \times 1$ vector. In particular, if the k -th covariate is of interest, \mathbf{a} can be written as $\mathbf{1}_k$, a $p \times 1$ vector with the k -th element 1 and rest 0. \otimes is the Kronecker product operator, and $\text{vec}(\cdot)$ is the vectorization of a matrix.

As an example, in our analysis of condylar shape data from 34 subjects (17 normal controls and 17 OA patients), we are interested in testing whether there is significant morphological difference between two groups (normal controls v.s. OA). We consider MFSDA (1) on the spatial x, y, and z coordinates of 1002 vertices on the condylar surface with $(y_{i1}(\mathbf{d}), y_{i2}(\mathbf{d}), y_{i3}(\mathbf{d}))^T = (\text{x coordinate}, \text{y coordinate}, \text{z coordinate})^T$, and $\mathbf{x}_i = (\text{intercept}, \text{group}, \text{age}, \text{gender}, \text{facial pain rate}, \text{the first principal component scores of all biological markers})^T$. In addition, the coefficient matrix \mathbf{C} in the hypothesis testing (2) can be written as $\mathbf{I}_3 \otimes \mathbf{1}_2^T$.

We introduce a local Wald-type test statistic $T_n(\mathbf{d})$ as follows:

$$T_n(\mathbf{d}) = \mathbf{r}(\mathbf{d})^T \left[\{\widehat{\Sigma}_\eta(\mathbf{d}, \mathbf{d})\}^{-1} \otimes \{[\mathbf{a}^T(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}]^{-1} \mathbf{a}^{\otimes 2}\} \right] \mathbf{r}(\mathbf{d}), \quad (3)$$

where $\mathbf{r}(\mathbf{d}) = \hat{\beta}(\mathbf{d}) - \text{Bias}(\hat{\beta}(\mathbf{d}))$, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$, and $\mathbf{a}^{\otimes 2} = \mathbf{a}\mathbf{a}^T$. The $\hat{\beta}(\mathbf{d})$ and $\hat{\Sigma}_\eta$ are estimates of the corresponding parameters, $\text{Bias}(\hat{\beta}(\mathbf{d}))$ is the bias term of $\hat{\beta}(\mathbf{d})$. Moreover, under the null hypothesis H_0 , the limiting distribution of $T_n(\mathbf{d})$ can be approximated by a weighted χ^2 distribution [Zhang and Chen, 2007].

To estimate all unknown parameters in model (1), we employ a weighted least squares (WLS) method based on the multivariate local polynomial kernel smoothing technique [Fan and Gijbels, 1996]. Let $K(\cdot)$ be a kernel function, and H be a bandwidth matrix with a simple diagonal form. We also denote that $K_{H,m}(\mathbf{d}) = |H|^{-1}K(H^{-1}(\mathbf{d}_m - \mathbf{d}))$ and $\mathbf{w}_H(\mathbf{d}_m - \mathbf{d}) = (1, (\mathbf{d}_m - \mathbf{d})^T H^{-1})^T$. For each j and a fixed bandwidth matrix H_β , the WLS estimator of $\beta_j(\mathbf{d})$ is given by

$$\hat{\beta}_j(\mathbf{d}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sum_{m=1}^{N_V} a_m(H_\beta, \mathbf{d}) \mathbf{y}_{\cdot,j}(\mathbf{d}_m), \quad (4)$$

where $a_m(H_\beta, \mathbf{d}) = \mathbf{e}^T [\sum_{m=1}^{N_V} K_{H_\beta,m}(\mathbf{d}) \{\mathbf{w}_{H_\beta}(\mathbf{d}_m - \mathbf{d})\}^{\otimes 2}]^{-1} K_{H_\beta,m}(\mathbf{d}) \mathbf{w}_{H_\beta}(\mathbf{d}_m - \mathbf{d})$, $\mathbf{e} = (1, \mathbf{0}_t^T)^T$, and $\mathbf{y}_{\cdot,j}(\mathbf{d}) = (y_{1,j}(\mathbf{d}), \dots, y_{n,j}(\mathbf{d}))^T$. Based on (4), for a fixed bandwidth matrix H_η , the WLS estimate of $\eta_{i,j}^{(g)}(\mathbf{d})$ is given by

$$\hat{\eta}_{i,j}^{(g)}(\mathbf{d}) = \sum_{m=1}^{N_V} a_m(H_\eta, \mathbf{d}) [y_{ij}(\mathbf{d}_m) - \mathbf{x}_i^T \hat{\beta}_j(\mathbf{d}_m)]. \quad (5)$$

Finally, we can estimate $\Sigma_\eta^{(g)}(\mathbf{d}, \mathbf{d}')$ by using the sample covariance function of $\hat{\eta}_i^{(g)}(\mathbf{d})$, denoted as $\hat{\Sigma}_\eta^{(g)}(\mathbf{d}, \mathbf{d}')$.

To select the optimal bandwidth in $\hat{\beta}_j(\mathbf{d})$ (or $\hat{\eta}_{i,j}(\mathbf{d})$), we use the generalized cross-validation score method [Zhu et al., 2012]. We standardize all covariates to have mean zero and standard deviation one; thus, we may choose a common bandwidth for all covariates. Moreover, following the arguments of Fan and Zhang [1999], a small bandwidth leads to a small value of $\text{Bias}(\hat{\beta}(\mathbf{d}))$, which can be dropped from the test statistics hereafter.

The global Wald-type statistic, denoted as T_n , is an integral of $T_n(\mathbf{d})\mu(\mathbf{d})$ with respect to $\mathbf{d} \in \mathcal{D}$; that is, $T(g) = \int_{\mathcal{D}} T_n(\mathbf{d})\mu(\mathbf{d})dL(\mathbf{d})$, where $L(\mathbf{d})$ is the Lebesgue measure. Selecting

different $\mu(\mathbf{d})$ allows us to introduce the prior information of specific regions of interest (ROIs). If there is no such prior information, then a uniform prior can be used. In this case, except for a constant scalar, T_n can be approximated by

$$T_n = \frac{Q_X^a}{N_V} \text{tr} \left(\left[\sum_{m=1}^{N_V} \mathbf{Y}_w(\mathbf{d}_m) \{ \hat{\Sigma}_\eta(\mathbf{d}_m) \}^{-1} \mathbf{Y}_w^T(\mathbf{d}_m) \right] \otimes [\mathbf{P}_X^a \mathbf{P}_X^{a^T}] \text{vec}(\mathbf{X}^T)^{\otimes 2} \right), \quad (6)$$

where $\mathbf{Y}_w(\mathbf{d}) = \sum_{m=1}^{N_V} a_m(H_\beta, \mathbf{d}) [\mathbf{y}_{1,\cdot}(\mathbf{d}_m), \dots, \mathbf{y}_{n,\cdot}(\mathbf{d}_m)]^T$, $Q_x^a = [\mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}]^{-1}$, and $\mathbf{P}_X^a = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}$. If the true region of interest, denoted by \mathcal{D}^* , is relatively large and its corresponding measurements are moderate, then the value of T_n should be relatively large. Thus, if the value of T_n is large, then morphological difference is far more likely to be found caused by covariates of interest.

We use the wild bootstrap method to approximate the null distribution of T_n under the assumption that the null hypothesis H_0 in (2) holds for all $\mathbf{d} \in \mathcal{D}$:

- Step 1. Fit model (1) under the null hypothesis H_0 , which yields $\hat{\underline{\beta}}^*(\mathbf{d})$, $\hat{\underline{\eta}}_i^*(\mathbf{d})$ and $\hat{\underline{\epsilon}}_i^*(\mathbf{d})$ for all i and \mathbf{d} .
- Step 2. Generate a random sample ν_i^b and $v_i^b(\mathbf{d}_m)$ from a $N(0, 1)$ generator for $i = 1, \dots, n$ and $m = 1, \dots, N_V$. B bootstrap samples are constructed as

$$\mathbf{y}_i^{(b)}(\mathbf{d}_m) = \mathbf{x}_i^T \hat{\underline{\beta}}^*(\mathbf{d}_m) + \nu_i^b \hat{\underline{\eta}}_i^*(\mathbf{d}_m) + v_i^b(\mathbf{d}_m) \hat{\underline{\epsilon}}_i^*(\mathbf{d}_m), \quad b = 1, \dots, B$$

for all i and $\mathbf{d}_m \in \mathcal{D}$.

- Step 3. For the b -th bootstrap samples, $b = 1, \dots, B$, calculate the global Wald-type statistic $T_n^{(b)}$ based on the formula in (6).
- Step 4. The p -value of T_n is calculated via using an approximation method [Huang et al., 2015]. In fact, T_n can be approximated by a χ^2 -type random variable $\alpha_1 \chi^2(\alpha_2) + \alpha_3$,

where α_1 , α_2 , and α_3 are respectively given by

$$\alpha_1 = \frac{\kappa_3(T)}{4\kappa_2(T)}, \quad \alpha_2 = \frac{8\kappa_2^3(T)}{\kappa_3^2(T)}, \quad \text{and} \quad \alpha_3 = \kappa_1(T) - \frac{2\kappa_2^2(T)}{\kappa_3(T)}, \quad (7)$$

where $\kappa_k(T)$, $k = 1, 2, 3$, are respectively the first three sample cumulants of $\{T_n^{(b)}\}_{b=1}^B$.

Finally, the p -value of T_n can be approximated by using $P(\chi^2(\alpha_2) \geq [T_n - \alpha_3]/\alpha_1)$.

References

- Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and Its Applications*. Chapman and Hall, London.
- Fan, J. and Zhang, W. (1999). Statistical estimation in varying coefficient models. *The Annals of Statistics*, 27(5):1491–1518.
- Huang, M., Nichols, T., Huang, C., Yang, Y., Lu, Z., Knickmeyer, R. C., Feng, Q., and Zhu, H. T. (2015). Fvgwas: fast voxelwise genome wide association analysis of large-scale imaging genetic data. *NeuroImage*, 118:613–627.
- Zhang, J. and Chen, J. (2007). Statistical inference for functional data. *The Annals of Statistics*, 35:1052–1079.
- Zhu, H., Fan, J., and Kong, L. (2014). Spatially varying coefficient model for neuroimaging data with jump discontinuities. *Journal of the American Statistical Association*, 109:977–990.
- Zhu, H., Kong, L., Li, R., Styner, M., Gerig, G., Lin, W., and Gilmore, J. H. (2011). Fadtts: functional analysis of diffusion tensor tract statistics. *NeuroImage*, 56:1412–1425.
- Zhu, H. T., Li, R. Z., and Kong, L. L. (2012). Multivariate varying coefficient model for functional responses. *Annals of Statistics*, 40:2634–2666.