

Extending People-like-Me Methods with Mahalanobis distance

A Personalized Predictive Longitudinal Study of Growth in Children

Randy (Xin) Jin

2023-06-20

1 Introduction

2 Methods

3 Results

4 Discussion

Objective

Modern clinical data analysis often deals with heterogeneous data

The traditional predictive modeling:

- Develops a model for global inference with all data
- Overlooks the diversity and heterogeneity

Personalized predictive methods

- Have the objective of addressing heterogeneity
- People-like-me method

People-like-me method:

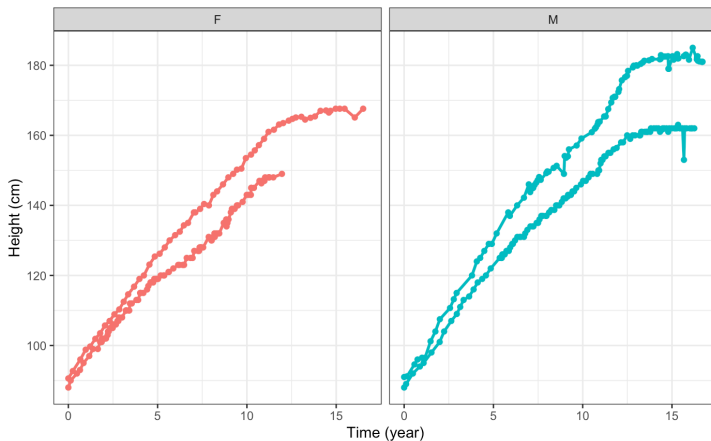
People-like-me:

- Uses a few but more similar trajectories
- Predicts based on those similar trajectories
- It can have higher predictive performance
- Previous work has focused on finding matching trajectories based on a single time point

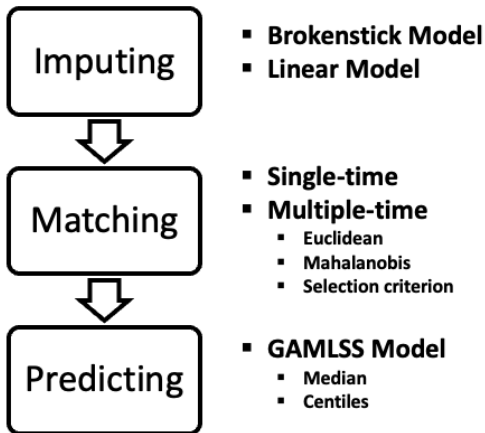
Our contribution is to:

- 1 extend matching step to using more time points (**anchor time**)
- 2 use alternative approaches to select matches

Example trajectories



Flowchart



Imputation

- 0 Split dataset into training and testing. For every individual in the testing set y_{test} , do the following steps.
- 1 Fit a brokenstick model with training set:

$$\mathbf{y}_{train} \sim BrokenStick(\mathbf{X}_{train})$$
 - brokenstick model is a particular type of linear mixed model
 - piece-wise cubic B-splines for both fixed and random effect
- 2 Obtain the predictions based on brokenstick model at fixed set of anchor time points $\boldsymbol{\tau} = \{\tau_1, \dots, \tau_k\}$: predict outcomes for every subjects i in training set: $\tilde{\mathbf{y}}_i = (\tilde{y}_{i,\tau_1}, \dots, \tilde{y}_{i,\tau_k})$

Imputation

- ③ Fit linear model with the predicted values from Step 2:

$$\tilde{\mathbf{y}} \sim LM_{\tau}(\mathbf{y}_{train,0}, \mathbf{X}_{train,0})$$

- This model includes only the demographic information and baseline outcome, at each time point independently
- ④ Obtain the predictions based on the linear model for every training subject $\bar{\mathbf{y}}_i$ and for the target individual $\bar{\mathbf{y}}_{test}$

Matching

- 5 Calculate the distance between the target (test) and every training subject: $D_{i,test} = dist(\bar{\mathbf{y}}_i, \bar{\mathbf{y}}_{test})$.

Two types of distance will be used

- Euclidean distance: $D_E = \|\mathbf{y}_1 - \mathbf{y}_2\| = \sqrt{\sum_{i=1}^k (y_{1i} - y_{2i})^2}$
- Mahalanobis distance: $D_M = \sqrt{(\mathbf{y}_1 - \mathbf{y}_2)^\top \Sigma^{-1} (\mathbf{y}_1 - \mathbf{y}_2)}$
 - where $\mathbf{y}_1 = (y_{11}, \dots, y_{1k})$ and $\mathbf{y}_2 = (y_{21}, \dots, y_{2k})$;
 - Σ is the variance-covariance matrix

Matching

- ⑥ Select the matching cohort as set \mathbf{A} from the training set, based on given distance in Step 5;
- The top κ smallest distance trajectories as matching set
- Mahalanobis distance criterion
 - $D_M^2 = (\mathbf{y}_1 - \mathbf{y}_2)^\top \Sigma^{-1} (\mathbf{y}_1 - \mathbf{y}_2)$, $D_M^2 \sim \chi_k^2$, where $df = k$ is the degrees of freedom
 - Then set \mathbf{A} may be chosen such that:

$$\mathbf{A} = \left\{ \mathbf{y} : (\mathbf{y} - \mathbf{y}^*)^\top \Sigma^{-1} (\mathbf{y} - \mathbf{y}^*) \leq \chi_{df=k, \alpha}^2 \right\}.$$

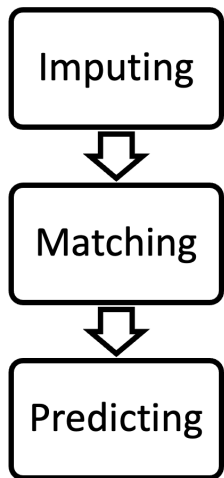
Prediction

- ⑦ Use the matching trajectories of $(\mathbf{y}_{sub}, \mathbf{X}_{sub}) \in \mathbf{setA}$ from Step 6 to fit a flexible model, specifically GAMLSS (Generalized Additive Model for Location Scale and Shape) model:
$$\hat{\mathbf{y}} \sim GAMLSS(\mathbf{X}_{sub})$$
- We use the predicted median as the personalized prediction and the centiles as the predictive interval for the target individual

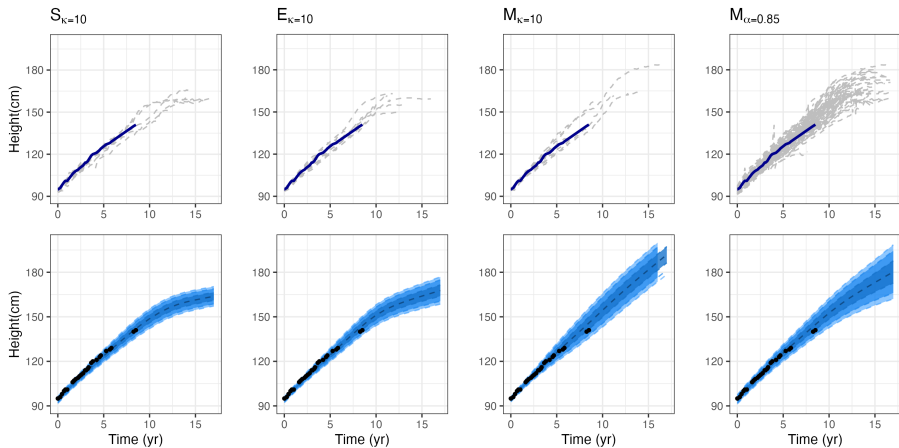
Application

The Early Pseudomonas Infection Control (EPIC) Observational Study is a prospective, multi-center, observational longitudinal study. EPIC study aims to investigate progression of cystic fibrosis outcomes in children

<i>Variable</i>	<i>Testing</i> , N = 457	<i>Training</i> , N = 913
<i>Female</i>	237 (52%)	456 (50%)
<i>Age₀</i> (<i>yr</i>)	3.14 (3.06, 3.23)	3.13 (3.05, 3.22)
<i>Age_{final}</i> (<i>yr</i>)	12.7 (10.3, 16.0)	12.8 (10.1, 16.0)
<i>Follow.up</i> (<i>yr</i>)	9.5 (7.2, 12.6)	9.7 (7.0, 12.7)
<i>Height₀</i> (<i>cm</i>)	94.0 (91.5, 96.7)	94.0 (91.4, 97.0)



- **Brokenstick Model**
 - fixed : piecewise cubic B-spline
 - random: piecewise cubic B-spline
 - knots = c(5, 10, 15)
- **Linear Model**
 - at each anchor time
 - covariates: baseline and demographics
- **Single-time:** $\kappa = 10, \tau = 10$
- **Euclidean:** $\kappa = 10$
- **Mahalanobis:** $\kappa = 10$
- **Mahalanobis:** $\alpha = 0.85$
- **GAMLSS Model**
 - mean: smoothing spline df = 3
 - sigma: smoothing spline df = 1



EPIC Data Analysis

	Anchor time	$S_{\kappa=10}^{t=10}$	$E_{\kappa=10}$	$M_{\kappa=10}$	$M_{\alpha=0.85}$
Bias	t(3, 6, 9, 12)	3.36	2.65	2.31	2.21
	t(4, 8, 12, 16)		2.22	2.29	2.20
RMSE	t(3, 6, 9, 12)	4.17	3.77	3.38	3.30
	t(4, 8, 12, 16)		3.43	3.49	3.29
50% CI	t(3, 6, 9, 12)	0.45	0.61	0.38	0.64
	t(4, 8, 12, 16)		0.53	0.47	0.64
80% CI	t(3, 6, 9, 12)	0.72	0.84	0.66	0.89
	t(4, 8, 12, 16)		0.79	0.74	0.89

Discussion

Summary:

- Improved people-like-me method for personalized prediction
- In the EPIC study, Mahalanobis distance ($\alpha = 0.85$) performed better than the other methods
- User-defined parameters: number of matches
- We also conducted a simulation study to assess performance of the methods
 - Mahalanobis distance ($\alpha = 0.85$) provide more accurate estimates, higher precision, and better coverage

Future work

- Imputation: use other prediction models
- Prediction Modeling: weighted regression and penalization
- Packaged functions for accessibility and application

Simulation

We performed a simulation study

- 1000 simulated datasets
- data generating mechanism based on linear mixed model
 - fixed effect: piece-wise cubic B-spline
 - random effect: piece-wise quadratic B-spline
- performing PLM approach

