

METRIC MULTI-DIMENSIONAL SCALING FOR LONGITUDINAL DATA EMBEDDINGS IN PHARMACOMETRICS

Anonymous authors

Paper under double-blind review

ABSTRACT

Longitudinal data in pharmacometrics typically involves multiple time-varying inputs and outputs for each subject in a population. Each subject can have a different number of observations at different time points, leading to irregular data structures that are difficult to analyze directly. Nonlinear mixed effects (NLME) models are the standard approach for modeling such data, but they can be computationally intensive and may not scale well with large datasets or complex models. In particular, for a large number of input covariates and output biomarkers and endpoints, the computational cost of fitting NLME models can become prohibitive. Some machine learning (ML) methods can be useful in eliminating useless covariates and biomarkers for a relatively low computational budget. Many ML models require fixed-size data as inputs and outputs. Such a tabular representation of a (usually) more complex data structure is commonly known as an embedding. In this work, we generate dissimilarity-preserving embeddings for longitudinal data commonly used in pharmacometrics. We use metric multi-dimensional scaling (MMDS) along with dynamic time warping (DTW) to generate fixed-size embeddings for each time-varying variable of each subject in a population. An experiment on a synthetic pharmacokinetic dataset shows that the proposed procedure can generate useful embeddings that preserve neighborhood structures. This has potential applications in covariate and biomarker elimination as well as model evaluation, to be investigated in future works.

1 INTRODUCTION

1.1 PHARMACOLOGY AND PHARMACOMETRICS

In pharmacology, drugs administered to the human body are considered to go through the ADME stages: absorption, distribution, metabolism, excretion. Pharmacometrics (Ette & Williams, 2006) is a branch of pharmacology dedicated to modelling these stages quantitatively through mathematical models. The data collected during clinical trials is longitudinal, consisting of time-varying observations for each subject in the study. This data is typically analyzed using ordinary differential equations (ODEs) in nonlinear mixed-effects (NLME) models (Owen & Fiedler-Kelly, 2014) to understand the drug’s behavior in the body and its effects over time across a population. The data includes doses given, covariates, drug concentration measurements and various biomarkers that change over time due to the drug’s effects.

1.2 BIOMARKERS AND COVARIATES

In clinical studies, there is typically a single primary endpoint that serves as the main outcome measured to assess the efficacy of a treatment. However, multiple secondary endpoints and biomarkers are often collected to provide additional insights into the drug’s effects and safety profile. Additionally, various covariates are recorded to account for individual differences among subjects that may influence the treatment response. Incorporating these additional variables into the analysis can enhance the predictive performance of pharmacometric models and enable a more precise characterization of drug behavior across diverse patient populations.

The covariates can be either baseline characteristics that do not change over time, such as age, weight, or genetic factors. Or they can be time-varying. On the other hand, biomarkers are typically

dynamic measurements that can vary over time in response to the drug treatment, such as blood pressure, heart rate, or specific protein levels. Often there may be 10s or 100s of biomarkers and covariates collected in a clinical trial, but only a few of them are truly informative for predicting the primary endpoint. Identifying and selecting the most relevant biomarkers and covariates is essential for building accurate pharmacometric models. Machine learning techniques can be used as a preliminary step to identify and eliminate uninformative variables, thereby reducing the final NLME model’s complexity.

1.3 MOTIVATION FOR EMBEDDINGS

The time-varying variables (covariates or biomarkers) often have a different number of observations per subject across variables and a different number of observations per variable across subjects. This irregular data structure poses challenges for traditional machine learning models that typically require fixed-size tabular data as inputs. While sequence models like recurrent neural networks (RNNs) (Elman, 1990; Werbos, 1990) and transformers (Vaswani et al., 2017) can handle variable-length sequences, they generally require more data and computational resources to train effectively. Given the limited data available in these pharmacometric studies, simpler machine learning models that operate on fixed-size tabular data may be preferred for their efficiency and interpretability.

2 METHODS

2.1 MULTI-DIMENSIONAL SCALING

Multi-dimensional scaling (MDS) (Borg & Groenen, 2005) is a set of techniques that can be used to create dissimilarity-preserving embeddings for each time-varying variable of each subject in a population. Common uses are dimensionality reduction and data visualization. Metric MDS (MMDS) is a variant of MDS that only requires a pairwise dissimilarity matrix between the points as an input. This makes it suitable for generating embeddings from longitudinal data, given a suitable dissimilarity function between subjects’ observations. The output of MMDS is a fixed-size embedding for each subject that preserves, as much as possible, the pairwise dissimilarities between the data objects for a given embedding dimension. MMDS formulates the pairwise dissimilarity preservation problem as a non-convex optimization problem, choosing the embeddings that locally minimize the sum of squares of the differences in dissimilarities between the original and embedding spaces. The MMDS formulation has the following definition, where d and δ are the dissimilarity functions in the input (longitudinal) and embedding spaces, respectively.

$$\min_Y \left(\frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n (d(x_i, x_j) - \delta(y_i, y_j))^2}{\sum_{i=1}^{n-1} \sum_{j=i+1}^n d(x_i, x_j)^2} \right)^{\frac{1}{2}}$$

For temporal variables, dynamic time warping (DTW) (Berndt & Clifford, 1994) can be used to compute pairwise dissimilarities between subjects based on their time-varying observations, even if the observation sequences do not have the same length. Notably, DTW is not a proper distance metric since it does not satisfy the triangle inequality. This means that MMDS cannot be guaranteed to find a perfect dissimilarity-preserving embedding. However, in practice, it can still produce useful embeddings that approximately preserve the pairwise dissimilarities.

We implemented MMDS and applied it to a synthetic pharmacokinetic (PK) (Owen & Fiedler-Kelly, 2014) dataset. A simple PK model was built and used to simulate observations for a virtual population of 34 subjects, each with 5 to 15 observations at random time points. More details about the specific model used and data simulated are provided in the supplementary materials. Each subject is represented as a multi-dimensional point defined by their sequence of observations over time. DTW was used to build a pairwise dissimilarity matrix. MMDS was applied to find the optimal embedding of each subject. The limited-memory BFGS optimizer (Nocedal & Wright, 2006) was used for optimizing the embeddings.

2.2 EMBEDDING EVALUATION

The embeddings were evaluated using the final MMDS loss objective as well as an experiment measuring how well neighborhood structures were preserved in the embedding space compared to the original space. To better understand the practical utility of the embeddings, we evaluated how well neighborhood structures were preserved in the embedding space compared to the original space. In particular, we found the k -nearest neighbors closest to each subject both in the original input space (using DTW dissimilarities) and in the embedding space (using Euclidean distances). This resulted in 2 sets of k neighbors for each subject. The 2 sets of neighbors were then compared for overlap. Additionally, we calculated the average DTW dissimilarity between each subject and its closest neighbors for each set of neighbors. The average dissimilarity was always calculated in the original space. The specific steps are laid out in Algorithm 1. A value of $k = 5$ was used in the experiment.

Algorithm 1 Embedding evaluation experiment

```

1: Set number of neighbors  $k$ 
2: Simulate observations from NLME model
3: Determine pairwise dissimilarities of observations using DTW
4: for embedding_dimension = 1 to max_dimension_of_embeddings do
5:   Initialize embeddings from standard normal distribution
6:   Scale embeddings by average pairwise dissimilarity and embedding dimension
7:   Optimize embeddings
8:   Calculate pairwise distances in embedding space
9:   for  $n = 1$  to number_of_subjects do
10:    Find  $k$  closest neighbors in original space
11:    Find  $k$  closest neighbors in embedding space
12:    Calculate, in original space, average DTW dissimilarity to embedding neighborhood
13:    Calculate, in original space, average DTW dissimilarity to original neighborhood
14:    return Size of intersection of both sets and the average dissimilarities
15:   end for
16: end for
17: Percentages are  $100 * (\text{sizes of intersections})/k$ 

```

3 RESULTS

Figure 1 shows the percentage of subjects that have 4 or 5 overlapping neighbors (out of $k = 5$) in their 2 sets of neighbors as a function of the embedding dimension. As can be seen in the top plot, a considerable overlap between the neighborhoods in both spaces was achieved. This indicates good structure preservation in the embedding procedure.

Figure 2 shows a scatter plot of the average dissimilarity from each subject to its neighbors defined in each space. Each subject is a point with the average dissimilarities representing the coordinates. A 45° is also included as a reference. Points lying on the line represent subjects that are equally distant to both sets of neighbors. As can be seen, good but not perfect agreement across spaces was achieved for embedding dimensions of 3 and beyond.

4 CONCLUSION AND FUTURE WORK

We used MMDS and DTW on longitudinal PK data to obtain tabular embeddings. An experiment was performed to validate that the neighborhood structure is preserved in the embedding space. The results indicate that the proposed procedure can generate useful dissimilarity-preserving embeddings for longitudinal data. In the future, we plan to apply the proposed embedding procedure along with tabular ML models to perform covariate and biomarker filtering. The goal is to identify uninformative variables that can be dropped before fitting a final NLME model. Another avenue of research is to use the embeddings for model evaluation by comparing the distributions of embeddings from real data and model simulations, similar to the Fréchet Inception Distance used in generative modelling

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

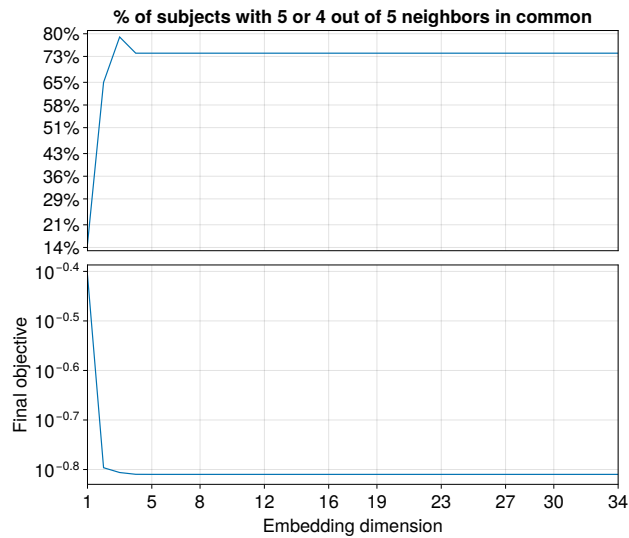


Figure 1: Top plot shows, for each embedding dimensionality, the percentage of overlap between neighborhoods in both spaces. Bottom plot shows final optimization objective values as a function of the embedding dimension.

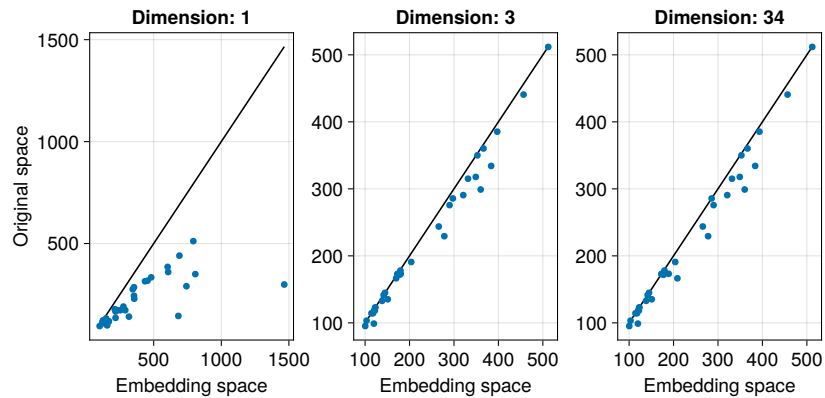


Figure 2: Average DTW dissimilarity to neighbors. Each scatter point represents a subject. The X and Y coordinates are a subject’s average DTW dissimilarity to its neighbors in the embedding space and original space, respectively. Proximity to the reference 45° line indicate structural agreement between spaces.

(Heusel et al., 2017). This could provide a new way to assess how well an NLME model captures the underlying data distribution.

REFERENCES

- Donald J. Berndt and James Clifford. Using dynamic time warping to find patterns in time series. *Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining*, pp. 359–370, 1994.
- Ingwer Borg and Patrick J. F. Groenen. *Modern Multidimensional Scaling: Theory and Applications*. Springer, 2005. doi: <https://doi.org/10.1007/0-387-28981-X>.
- Jeffrey L. Elman. Finding structure in time. *Cognitive Science*, 14(2):179–211, 1990. ISSN 0364-0213. doi: [https://doi.org/10.1016/0364-0213\(90\)90002-E](https://doi.org/10.1016/0364-0213(90)90002-E).

- 216 Ene I. Ette and Paul J. Williams. *Pharmacometrics: The Science of Quantitative Pharmacology*.
 217 John Wiley & Sons, Inc., 2006. doi: <https://doi.org/10.1002/0470087978>.
 218
- 219 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and
 220 Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilib-
 221 rium. In *Proceedings of the 31st Conference on Neural Information Processing Systems (NeurIPS*
 222 *2017)*, 2017.
- 223 Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006. doi:
 224 <https://doi.org/10.1007/978-0-387-40065-5>.
 225
- 226 Joel S. Owen and Jill Fiedler-Kelly. *Introduction to Population Pharmacokinetic / Pharmacody-*
 227 *namic Analysis with Nonlinear Mixed Effects Models*. John Wiley & Sons, Inc., 2014.
- 228 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 229 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Infor-*
 230 *mation Processing Systems*, volume 30, 2017.
- 231 P.J. Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*,
 232 78(10):1550–1560, 1990. doi: 10.1109/5.58337.
 233

234 A APPENDIX

235 A.1 NONLINEAR MIXED EFFECTS MODELS NOTATION

236
 237 Let each subject $i \in \{1, 2, \dots, n\}$ be independently measured at time points $t_{i,1}, t_{i,2}, \dots, t_{i,m_i}$,
 238 leading to the vector of observations $Y = \{Y_{i,j}, \text{ for } i \in 1 \dots n \text{ and } j \in 1 \dots m_i\}$. For simple
 239 Gaussian response variables, the response over time is given by $Y = f(x, \beta) + g(x, \beta, \sigma) \odot \epsilon$, where
 240

- 241 • x consists of time t and subject-specific covariates (time-varying or baseline), including
 242 baseline covariates b as a component,
- 243 • β is a vector of estimable parameters, some (random effects) varying across subjects and
 244 others not (fixed effects),
- 245 • $f(x, \beta)$ is a function whose output has the shape as Y defining the structural model, often
 246 defined by a system of ODEs,
- 247 • ϵ is a vector of independently sampled errors of the same shape as Y , following a standard
 248 normal distribution $\epsilon_{i,j} \sim \mathcal{N}(0, 1)$,
- 249 • σ is a vector of residual error parameters, shared for all subjects,
- 250 • $g(x, \beta, \sigma)$ is a function whose output has the same shape as Y defining the residual error
 251 model, and
- 252 • \odot is the element-wise product operator.

253 Let the component of the parameters β that is varying across subjects be $\eta_i \sim \mathcal{N}(h(\theta, b_i), \Omega)$ (one
 254 vector for each subject i), where

- 255 • θ is a vector of estimable parameters, shared for all subjects,
- 256 • b_i is the set of baseline covariates for subject i ,
- 257 • $h(\theta, b_i)$ is a function defining the relationship between the baseline covariates for a partic-
 258 ular subject and the individual parameters' mean, and
- 259 • Ω is another set of parameters representing the variance-covariance matrix of η_i , shared for
 260 all subjects.

261 For simplicity, let the remaining part of β that is not varying between subjects be part of θ . In other
 262 words, β can be derived from θ and η_i , for all subjects i , and the only estimable parameters of the
 263 model are:

- 264 1. Population-level parameters θ , Ω , and σ , and

270 2. Subject-level parameters η_i for each subject i .

271
272 In practice, the PK parameters β are often positive-valued so $f(x, \beta)$ will often involve exponen-
273 tiating the corresponding components of β . For example, the clearance rate (the rate at which the
274 body metabolizes the drug) of subject i can be defined as $CL_i = \exp(\eta_{i,1})$ where $\eta_{i,1}$ is the first
275 component of η_i with mean $h(\theta, b_i)_1$. It is also common to assume that the η_i s have a mean of zero,
276 adding their mean value as part of $f(x, \beta)$ instead, e.g.:

$$277 \eta_i \sim \mathcal{N}(0, \Omega)$$

$$278 CL_i = \exp(h(\theta, b_i)_1 + \eta_{i,1}) = \exp(h(\theta, b_i)_1) \cdot \exp(\eta_{i,1})$$

279
280 If there are no baseline covariate effects and $h(\theta, b_i)_1 = \theta_1$, this simplifies to $CL_i = \theta_{CL} \cdot \exp(\eta_{i,1})$
281 where $\theta_{CL} = \exp(\theta_1)$. Since θ_1 is an estimable parameter, we can re-parameterize it as $\theta_{CL} > 0$
282 directly. In this case, θ_{CL} is the population-level clearance parameter, while $\eta_{i,1}$ represents subject
283 i 's log-scale variation from that population value.

284 The residual error model defined by $g(x, \beta, \sigma)$ can take different forms, e.g.:

- 285
286 1. Additive, where the residual error is constant across all time points and subjects,

$$287 g(x, \beta, \sigma)_{i,j} = \sigma_1$$

- 288
289 2. Proportional, where the residual error scales with the predicted value $f(x, \beta)$,

$$290 g(x, \beta, \sigma)_{i,j} = \sigma_1 \cdot f(x, \beta)_{i,j}$$

- 291
292 3. Combined, where both additive and proportional components are present

$$293 g(x, \beta, \sigma)_{i,j} = \sqrt{\sigma_1^2 + (\sigma_2 \cdot f(x, \beta)_{i,j})^2}$$

294
295
296 A.2 PK MODEL USED IN SIMULATION

297
298 For the experiments presented, PK data was simulated using a simple two-compartment (Central and
299 Peripheral) PK NLME model with combined residual errors. The only covariates in the model were
300 time and the drug's dose. The latent encoding of random effects η are considered independent and
301 are sampled from a multivariate normal distribution with zero mean and diagonal covariance matrix
302 Ω with non-zero elements $[0.25, 0.25, 0.25, 0.25, 0.49]$. The individual parameters are defined below
303 where:

- 304
305 • CL is the clearance rate,
306 • Vc is the hypothetical volume of the Central compartment,
307 • Q is the intercompartmental clearance which is the rate of transfer of the drug from the
308 Central to the Peripheral compartment,
309 • Vp is the hypothetical volume of the Peripheral compartment;
310 • Ka is the absorption rate constant from Depot to Central.

$$311$$

$$312 CL_i = CL \cdot \exp(\eta_{i,1})$$

$$313 Vc_i = Vc \cdot \exp(\eta_{i,2})$$

$$314 Q_i = Q \cdot \exp(\eta_{i,3})$$

$$315 Vp_i = Vp \cdot \exp(\eta_{i,4})$$

$$316 Ka_i = Ka \cdot \exp(\eta_{i,5})$$

317
318 The vector of population-level parameters θ is thus $[CL, Vc, Q, Vp, Ka]$. The specific values of the
319 parameters used in the simulation are included in Table 1.

320
321 In the PK model used, the Depot refers to the a compartment associated with the gut for orally-
322 administered drugs. The Central compartment is linked to richly perfused organs and blood plasma.
323 And the Peripheral compartment represents slowly perfused organs, like fat and skin. A compart-
ment name with a prime refers to the rate of change of drug concentration therein.

Table 1: Values of fixed effects θ used in simulation.

Parameter in θ	Value
Clearance (CL)	25
Central volume (V_c)	50
Intercompartmental clearance (Q)	31
Peripheral volume (V_p)	48
Absorption constant (K_a)	0.5

In this simulation, the drug is administered orally so the dose is given into the Depot compartment as a bolus at time 0. The drug amounts (not concentration) over time in each compartment are defined by the following system of ODEs:

$$\begin{aligned} \text{Depot}' &= -K_a \cdot \text{Depot} \\ \text{Central}' &= K_a \cdot \text{Depot} - \frac{\text{CL}}{V_c} \cdot \text{Central} - \frac{Q}{V_c} \cdot \text{Central} + \frac{Q}{V_p} \cdot \text{Peripheral} \\ \text{Peripheral}' &= \frac{Q}{V_c} \cdot \text{Central} - \frac{Q}{V_p} \cdot \text{Peripheral} \end{aligned}$$

Lastly, Y represents the drug concentration measurements in the Central compartment. The residual error model $g(x, \beta, \sigma)$ is defined such that a combined residual error model is used. The additive and proportional components of the residual error parameters $\sigma = [\sigma_{\text{add}}, \sigma_{\text{prop}}]$ used in the simulation were both 0.05.

$$\begin{aligned} C_p &= \frac{\text{Central}}{V_c} \\ Y &\sim \mathcal{N}(C_p, (C_p \cdot \sigma_{\text{prop}})^2 + \sigma_{\text{add}}^2) \end{aligned}$$

Figure 3 shows the simulated trends of drug concentrations for all subjects.

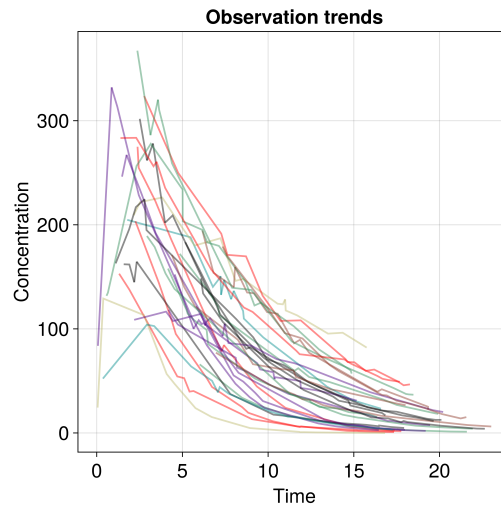


Figure 3: Concentration trends for all subjects in simulated population.

A.3 LARGE POPULATION

It's difficult to retain the structure across spaces when simulating data with larger populations (here, 75 subjects). This is shown in Figure 4 by the lower percentage of subjects with large intersections of neighbor sets.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

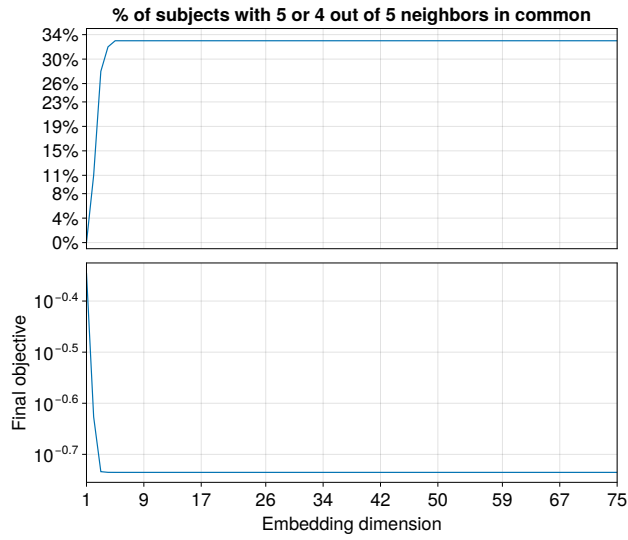


Figure 4: Results of neighborhood overlap experiment with larger population of 75 subjects.

However, as Figure 5 details, the distance experiment still presents good results.

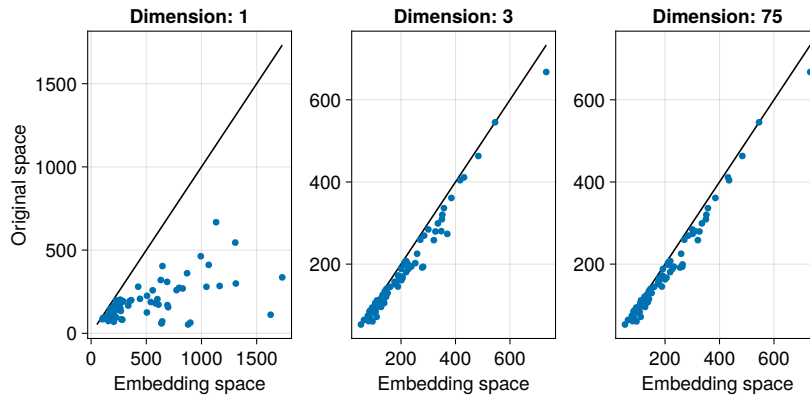


Figure 5: Results of neighborhood distances experiment with larger population of 75 subjects.