# Three-dimensional Radial Visualization of High-dimensional Continuous or Discrete Datasets

Fan Dai, Yifan Zhu and Ranjan Maitra

Department of Statistics
Iowa State University
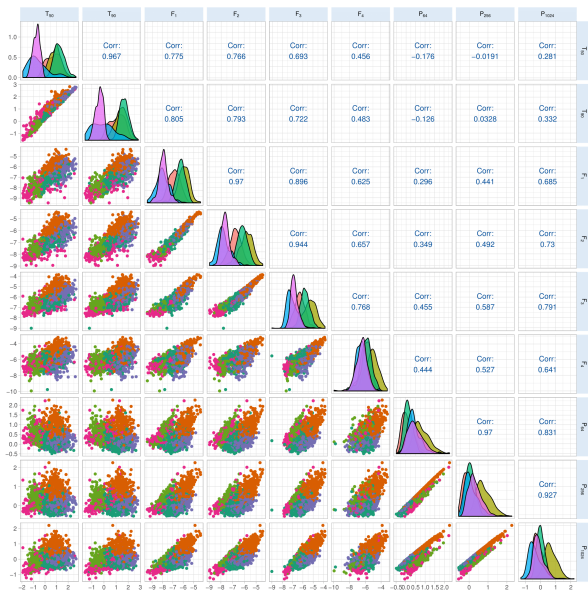{fd43,yifanzhu,maitra}@iastate.edu

# Motivation

- Multivariate datasets
  - agriculture, engineering, genetics, social science. . .
- Complex data structure
  - datasets with many discrete, skewed or correlated features
    - image, voice, surveys. . .
    - need advanced methods for analysis and summaries
- Display distinct groups while also inherent variability

# Example: Gamma Ray Bursts (GRBs)

- Extremely energetic explosions observed in distant galaxies.
  - data from NASA's Burst and Transient Source Experiment
  - 1,599 GRBs with complete information on 9 parameters
    - time for % flux to arrive, peak fluxes in different channels, time-integrated fluences over time-points
- Nine heavily-skewed "parameters" or attributes
  - use of logarithms to reduce skewness
- astrophysics community argued long over 2 or 3 types
  - analysis based on summary exclusion of some heavily-correlated attributes
  - recent analysis shows all 9 features important for clustering
    - actually 5 ellipsoidal groups, not 2 or 3
- smaller-dimensional 9D example used as a test case

- Visualization tools for continuous multivariate data
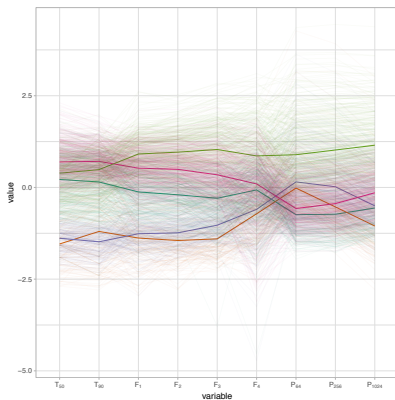  - pairwise scatter plots

# Pairwise Scatterplots: Gamma Ray Bursts

# Background and Current Work

- Visualization tools for continuous multivariate data
    - pairwise scatter plots
        - limited in providing multivariate assessments
    - parallel coordinates plot (*Inselberg '85, Wegman '90*)

# Parallel Coordinate Plots: Gamma Ray Bursts



- Represent multidimensional data using lines.

  - vertical line represents each dimension or attribute.
  - $p - 1$ lines connected at appropriate scaled dimensional value represent each observation
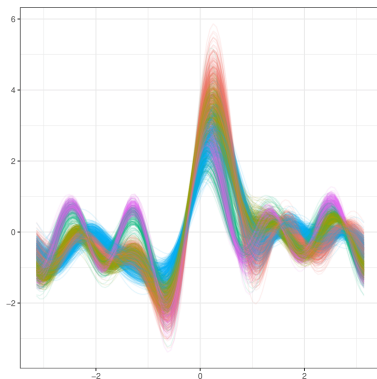  - polar version provided by star plot

# Background and Current Work

- Many approaches to display continuous multivariate data

    - pairwise scatter plots
        - limited in providing multivariate assessments
    - parallel coordinates plot (*Inselberg '85, Wegman '90*)
        - placement order matters, unclear for large $n$, $p$
        - hard to identify groups/patterns with even moderate $n$.
    - Andrews' curves represent each observation via trigonometric series

# Andrews' Curves: Gamma Ray Bursts

- Plot each $\boldsymbol{X} = (X_1, X_2, \ldots, X_p)$ as a curve:

$$f(t) = x_1 + x_2 \sin t + x_3 \cos t + x_4 \sin 2t + x_5 \cos 2t + \ldots, \qquad t \in [-\pi, \pi]$$



- Entire curve displays one observation

# Background and Current Work

- Many approaches to display continuous multivariate data

    - pairwise scatter plots
        - limited in providing multivariate assessments
    - parallel coordinates plot (*Inselberg '85, Wegman '90*)
        - placement order matters, unclear for large $n, p$
        - polar version provided by star plot
    - Andrews' curves
        - order in which coordinate enters series important
        - very computationally intensive for larger $p$
    - Star coordinates plot
      represents coordinate axes as equi-angled rays extending from center
        - order matters, optimized (*van Long & Linsen '11*)

- Use springs to display observation (radial visualization)

# Two-dimensional radial visualization (RadViz2D)

- Uses Hooke's law to project data onto unit circle
  - place *p* springs (anchor points) on the rim
    - pull each spring by value relative to coordinate from center
    - observations w/ similar relative values in all attributes end up closer to center, others are closer to the anchor points
  - order of placement of springs affects display
    - refinements to improve RadViz2D exist (see later)

# RadViz2D Illustration

$\boldsymbol{X} = (X_1, X_2, X_3, X_4, X_5) = (0.7, 0.5, 0.3, 0.2, 0.7)$

- Maps $\boldsymbol{X} \in \mathbb{R}^p$ to 2D point $\boldsymbol{\Psi}^\bullet(\boldsymbol{X}; \boldsymbol{U}) = \boldsymbol{UX}/\boldsymbol{1}_p'\boldsymbol{X}$:
  $\boldsymbol{U}$ projection matrix, columns (anchor points) on $\mathbb{S}^1$

# Two-dimensional radial visualization (RadViz2D)

- Uses Hooke's law to project data onto unit circle

    - place *p* springs (anchor points) on the rim

        - pull each spring by value relative to coordinate from center
        - observations w/ similar relative values in all attributes end up closer to center, others are closer to the anchor points

    - order of placement of springs affects display

        - refinements to improve RadViz2D exist (see later)

- Effective for sparse data, in evaluating distinct groups

    - Nonlinear map distorts, affects interpretability
    - High-dimensional observations more difficult to visualize

- Can fully 3D extension improve performance?

    - Viz3D provides third dimension, constant for all observations (*Artero & de Oliveira, '04*)

# Generalizing Radial Visualization

- Allow anchor points in $U$ on $\mathbb{S}^q$, $q > 1$, not necessarily equi-spaced

  - $p$ springs at $u_1, u_2, \ldots, u_p \in \mathbb{S}^q$, with spring constants $X_1, X_2, \ldots, X_p$.
  - equilibrium point $Y \in \mathbb{R}^{q+1}$ of system satisfies

    $$\sum_{j=1}^{p} X_j(Y - u_j) = 0,$$

    - $Y = \Psi(X; U) = UX/1'_p X$ solves the system.
  - is line-, point-ordering- and convexity-invariant.
  - scaling every coordinate to be in [0,1] allows for $Y \in \mathbb{S}^q$.

# Placement of Anchor Points

- Suppose: coordinates of $X$ are uncorrelated.
- For $X_1, X_2 \in \mathbb{R}^p$, let $Y_i = \Psi(X_i; U), i = 1, 2$.

  - Euclidean distance between $Y_1$ and $Y_2$ is

  $$\| Y_1 - Y_2 \|^2 = \left( \frac{X_1}{1_p' X_1} - \frac{X_2}{1_p' X_2} \right)' U' U \left( \frac{X_1}{1_p' X_1} - \frac{X_2}{1_p' X_2} \right),$$

    - $X_i$, $X_j$ very dissimilar, with perfect negative correlation, should be placed as far away as possible (in opposite directions) in our radial visualization.

- However, $\| Y_i - Y_j \|^2 \to 0$ as $\langle u_i, u_j \rangle \to 0$.

  - may create artificial visual correlation between $i$th and $j$th coordinates if $\langle u_i, u_j \rangle \to 0 < \pi/2$.
  - need $u_j$s far from the other as possible; so evenly distributed.
  - $\mathbb{S}^q$: for larger $q$, can get larger angles between $u_j$s

- Also place positively correlated coordinates close together

  - $q > 1$ has advantage in placing multiple coordinates together
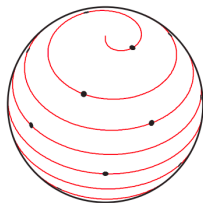
# Three-dimensional Radial Visualization

- $q = 2$ in our generalization yields $\mathrm{RadViz3D}$:

  - equi-spaced anchor points for 5 Platonic solids, $p = 4, 6, 8, 12, 20$.

    - closely related to Thomson problem in traditional molecular quantum chemistry (Atiyah & Sutcliffe '03).

  - for other $p$, approximate through Fibonacci grid, $j$th anchor point:

  $$u_{j1} = \cos(2\pi j \varphi^{-1})\sqrt{1 - u_{j3}^2},$$
  $$u_{j2} = \sin(2\pi j \varphi^{-1})\sqrt{1 - u_{j3}^2},$$
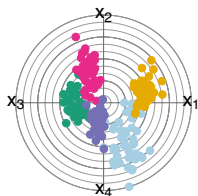  $$u_{j3} = \frac{2j - 1}{p} - 1,$$

  where $\varphi = (1 + \sqrt{5})/2$ is the golden ratio.

  

  (González '10)

    - distributes anchor points along generative spiral on $\boldsymbol{S}^2$, with consecutive points as separated as possible, satisfies "well-separation" property (Saff & Kuijlaars '97).
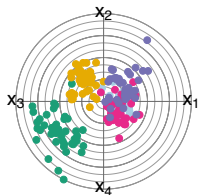
# 4D Examples simulated via *MixSim* package in R



RadViz2D, $\ddot{\omega} = 10^{-3}$

Viz3D            RadViz3D



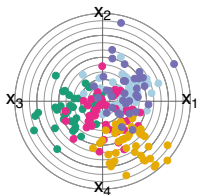RadViz2D, $\ddot{\omega} = 10^{-2}$

# Higher-dimensional Datasets

- Display *p* anchor points infeasible, even for moderate *p*

  - placement of equally-spaced anchor points built on not inducing spurious positive correlations in display

    - with increasing *p*, harder to guarantee such outcome

- Project high-dimensional data to uncorrelated coordinates but preserve distinctiveness and variability in groups

  - Principal Components finds mutually orthogonal projections summarizing proportion of total variance, but does not account for groups.

# Maximum-Ratio Projection (MRP)

- **Step 1:** Obtain PCs (orthogonal $V_g$) for each group

  - Find orthogonal $W$ closest to all $V_g$

    - Project $X$ with $W$ and then obtain MRP

- **Step 2:** Obtain uncorrelated projections that maximize between-group sums of squares and cross products (SSCP) relative to the total SSCP.

  - Let $T$, $W$ be (p.d.) total & between-group corrected SSCP.

    - $\hat{v}_j = T^{-\frac{1}{2}} \hat{w}_j / \| T^{-\frac{1}{2}} \hat{w}_j \|, j = 1, 2, \ldots, k$, $\hat{w}_j, j = 1, 2, \ldots, k$ are, in decreasing order, the $k$ largest eigenvalues of $T^{-1/2} B T^{-1/2}$.
    - $k \leq G - 1$, chosen by scree plot/quality of display
    - $G \leq 4$ needs $4 - G + 1$ more projections w/ null contribution
    - needs p.d. $T$, does not hold if $p > \min n_g$

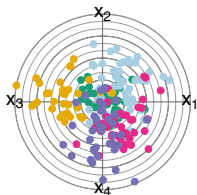- MRP maximizes separation between groups (in projected space) relative to total variability.

# 500D Examples



RadViz2D, $\ddot{\omega} = 10^{-3}$
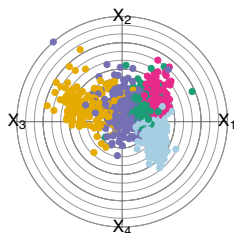
Viz3D                     RadViz3D



RadViz2D, $\ddot{\omega} = 10^{-2}$

# Datasets with Skewed Attributes

- Consider a r.v. $X$ with CDF $F_X(x)$.

    - $F_X(X) \sim U(0,1) \Rightarrow Y = \Phi^{-1}[F_X(X)] \sim N(0,1)$.
        - call the above (classical) Gaussianized Distributional Transform (CGDT)
        - marginal application of CGDT specifies distribution on $X$ with desired marginal and correlation structure.

- CGDT standardizing transform, more stringent than usual affine 0-mean, unit-variance inducing transform

    - CGDT matches all marginal quantiles to N(0,1)

    - Apply to skewed datasets or with unclear marginals

- Before applying MRP and RadViz3D

# Applications: Gamma Ray Bursts Dataset



RadViz2D

Viz3D                    RadViz3D

Groups ● 1 ● 2 ● 3 ● 4 ● 5

- Heavily skewed attributes, so CGDT appropriate

- Results indicate 5 overlapping clusters

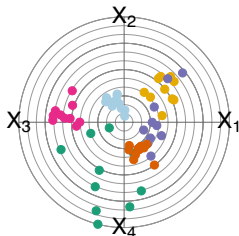  - some suggestion of 2, 3 super-types of GRBs

# Applications: Face Recognition



- $112 \times 92$-images of 6/40 faces at 10 light angles/conditions.
- $(20 \times 14)$ DWT2 (LL band) of wavelet-transformed images with 280 features (Jadhav & Holambe, 2009)

# Applications: Face Recognition



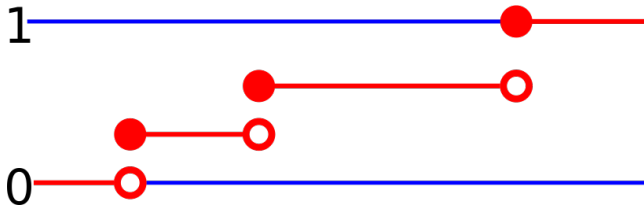Persons ● A ● B ● C ● D ● E ● F



RadViz2D

Viz3D                    RadViz3D

- marginals unclear: use CGDT

- RadViz3D clarifies all 6 people the best

# Datasets with Discrete Attributes

- For discrete-valued variable $X$, CDF $F_X(X) \not\sim U(0, 1)$ because of discreteness.

  - CGDT currently not applicable



- Note that the CDF is only right continuous

- Solution proposed by Rüschendorf (2013) via the generalized distributional transform

# Generalized Distributional Transform (GDT)

## Definition

Let $X$ be a real-valued RV with CDF $F_X(\cdot)$ and let $V \sim U(0, 1)$ be a RV independent of $X$. The generalized distributional transform of $X$ is $U = F(X, V)$ where $F(x, \lambda) \doteq P(X < x) + \lambda P(X = x) = F_X(x-) + \lambda[F_X(x) - F_X(x-)]$ is the generalized CDF of $X$.
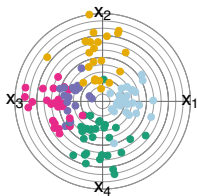
## Theorem

*Let $U = F(X, V)$ be the generalized distributional transform of $X$. Then*

$$U \sim \text{Uniform}(0, 1) \text{ and } X = F_X^{-1}(U) \text{ a.s.}$$

*where $F^{-1}(t) = \inf\{x \in \mathbb{R} : F_X(x) \geq t\}$ is the generalized inverse, or the quantile transform, of $F_X(\cdot)$.*

- Use $F(X, V)$ in place of $F_X(X)$, calculate GDT as before
  - use of GDT on non-discriminating coordinate can spuriously bestow it hyper-importance
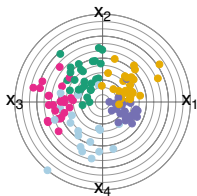    - suggest ANOVA test on each GDT-ed coordinate, control FDR

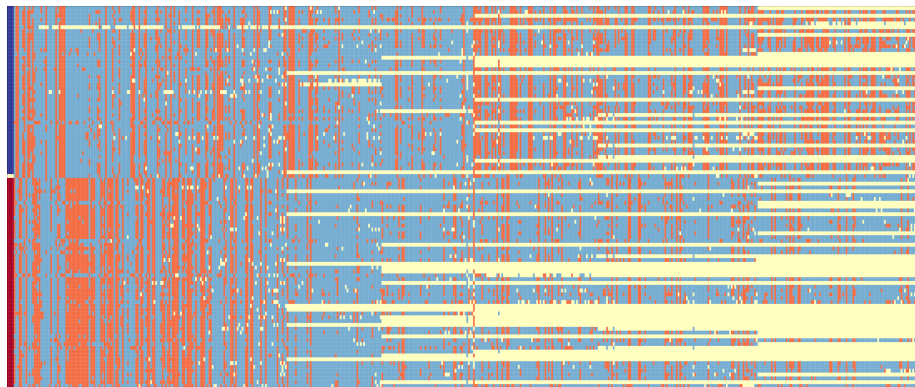# Illustration: Simulated Binary Datasets



RadViz2D, low
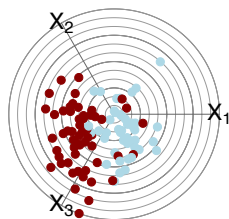clustering complexity

Viz3D            RadViz3D



RadViz2D, high
clustering complexity

# Applications: Senate Voting Records



- 108th US Congress (2005-06) had 542 (Y/N/NR) Senate votes

  - 55 Republicans, 44 Democrats, 1 (D-caucus) Independent (VT)
    (Banerjee *et al*, 2008)

    - combine N/NR to get dataset of binary attributes

# Applications: Senate Voting Records



RadViz2D

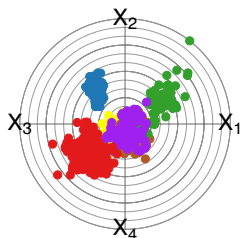Viz3D       RadViz3D

● Democratic ● Republican

- $G = 2$ so only 1 MRP with postive eigenvalue

  - spring $X_1$ pulls members of one party towards itself more
  - $X_2, X_3, X_4$ pull senators from both parties with equally (non-discriminating) force

# Applications: Handwritten Indic Scripts



(Map Acknowledgment: Surveyor-General of India)

- Handwritten scripts from Bangla (east), Gujarati (west), Gurmukhi (north), Kannada and Malayalam (southern states of Karnataka and Kerala), Urdu (Persian script), with 116 mixed features (Obaidullah et al, 2017).

# Applications: Handwritten Indic Scripts



RadViz2D
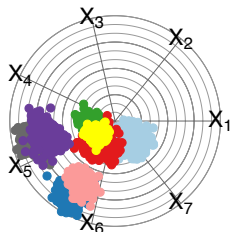
Viz3D                    RadViz3D

● Bangla ● Gujarati ● Gurmukhi ● Kannada ● Malayalam ● Urdu

- Viz3D (lesser extent RadViz2D) separates Urdu, Kannada and Gujarati, not the other 3 languages

- RadViz3D best in classifying all the 6 scripts
  - also points to difficulty of problem

# Applications: *RNA* Sequences

- Gene expression levels, in FPKM, of RNA sequences from 13 human organs.

  - focus on 8 largest (in terms of the sample size) organs

    - esophagus (659), colon (339), thyroid (318), lung (313), breast (212), stomach (159), liver (115) and prostate (106)

  - $p$=20242 discrete features

    - some have many discrete values, essentially continuous

  - dataset of mixed attributes.

- Display for distinctiveness of samples from each organ

# Applications: RNA Sequences



RadViz2D

Viz3D                    RadViz3D

● Breast ● Colon ● Esophagus ● Liver ● Lung ● Prostate ● Stomach ● Thyrioid

- RadViz2D, Viz3D poorer at separating organs

- RadViz3D indicates clear separation between organs

  - colon and stomach have some marginal overlap.

# Conclusions and Further Work

- Visualization tool for HD datasets

  - RadViz3D for more comprehensive display of grouped data

  - MRP, GDT for discrete, mixed, skewed variates

  - displays distinct groups more accurately

  - R package `https://github.com/fanne-stat/radviz3d`

  - manuscript `https://arxiv.org/abs/1904.06366`

- Number of issues merit further attention

  - MRP linear; non-linear projections better?

  - extend for categorical (non-binary) attributes

  - GDT/MRP with other tools for improved visualization