

# 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

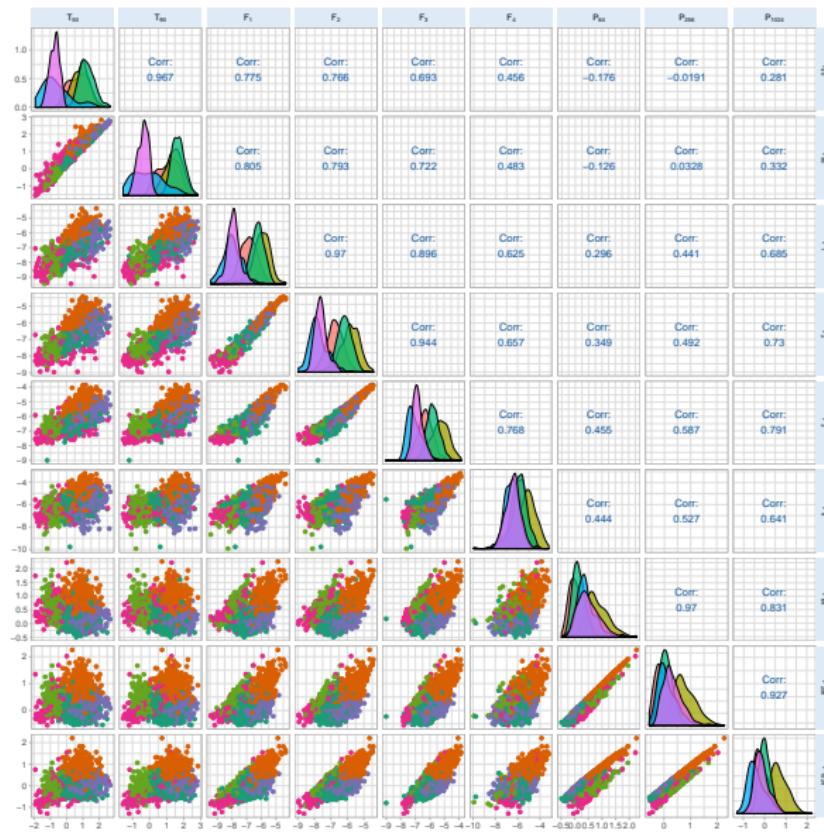
- Multivariate datasets
  - agriculture, engineering, genetics, social science...
- Complex data structure
  - datasets with many discrete, mixed or skewed features
    - image, voice, surveys...
    - need advanced methods for analysis and summaries
- Display distinct groups while also inherent variability

## Example: Gamma Ray Bursts Dataset

- Brightest electromagnetic events in space, believed to contain clues to the origin of the cosmos
  - 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

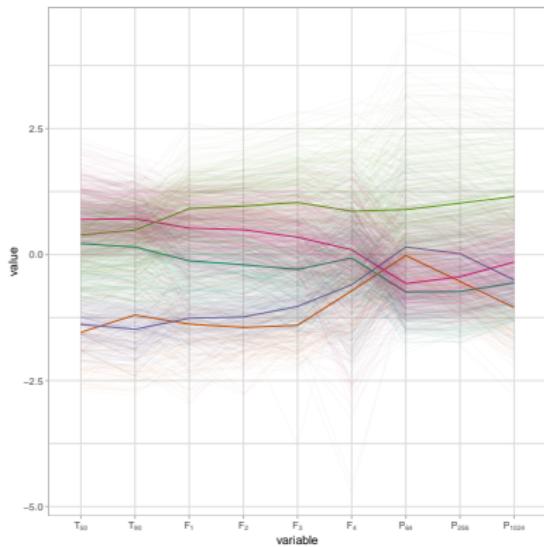
# Pairwise Scatterplots: Gamma Ray Bursts



# 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
- Not easy to identify groups/patterns with even moderate  $n$ .

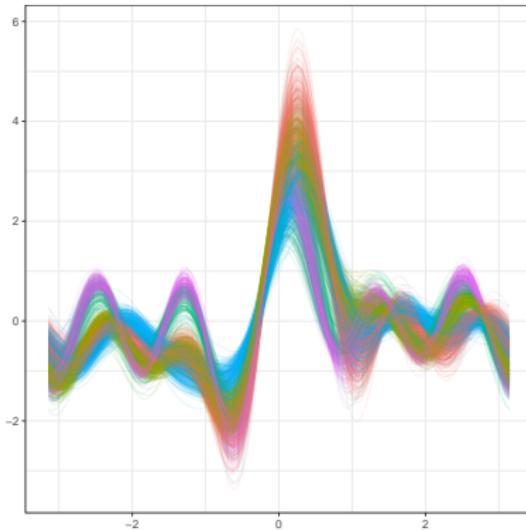
# Background and Current Work

- 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 represent each observation via trigonometric series

# Andrews' Curves: Gamma Ray Bursts

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

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



- 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 represent each observation via trigonometric series
    - order in which coordinate enters series important
    - very computationally intensive for larger  $p$
  - surveyplot represents each observed feature as line graph of length relative to size
  - 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)
- 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*)

## RadViz2D Illustration

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

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

# Generalizing Radial Visualization

- Allow anchor points in  $\mathbf{U}$  on  $\mathbb{S}^q$ ,  $q > 1$ , not necessarily equi-spaced (yet!)
  - physical interpretation, has  $p$  springs at  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p \in \mathbb{S}^q$ , with spring constants  $X_1, X_2, \dots, X_p$ .
  - equilibrium point  $\mathbf{Y}$  of system satisfies

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

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

# Placement of Anchor Points

- Suppose: coordinates of  $\mathbf{X}$  are uncorrelated.
- For  $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^p$ , let  $\mathbf{Y}_i = \Psi(\mathbf{X}_i; \mathbf{U}), i = 1, 2$ .

- Euclidean distance between  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$  is

$$\|\mathbf{Y}_1 - \mathbf{Y}_2\|^2 = \left( \frac{\mathbf{X}_1}{\mathbf{1}'_p \mathbf{X}_1} - \frac{\mathbf{X}_2}{\mathbf{1}'_p \mathbf{X}_2} \right)' \mathbf{U}' \mathbf{U} \left( \frac{\mathbf{X}_1}{\mathbf{1}'_p \mathbf{X}_1} - \frac{\mathbf{X}_2}{\mathbf{1}'_p \mathbf{X}_2} \right),$$

i.e. p.d. quadratic form; unit vector columns ( $\mathbf{u}_j$ ) of  $\mathbf{U}$

- $(i, j)$ th element of  $\mathbf{U}' \mathbf{U}$ :  $\mathbf{u}'_i \mathbf{u}_j = \cos\langle \mathbf{u}_i, \mathbf{u}_j \rangle$ ;  $\mathbf{u}'_i \mathbf{u}_i = 1$ .
- For  $\mathbf{X}_i = a_i \mathbf{e}_i, \mathbf{X}_j = a_j \mathbf{e}_j$ ,  $\|\mathbf{Y}_i - \mathbf{Y}_j\|^2 = 2 - 2 \cos\langle \mathbf{u}_i, \mathbf{u}_j \rangle$ .
- $\mathbf{X}_i, \mathbf{X}_j$  very dissimilar, with perfect negative correlation, so should be placed as far away as possible (in opposite directions) in our radial visualization.
- But  $\|\mathbf{Y}_i - \mathbf{Y}_j\|^2 \rightarrow 0$  as  $\langle \mathbf{u}_i, \mathbf{u}_j \rangle \rightarrow 0$ .
  - may create artificial visual correlation between  $i$ th and  $j$ th coordinates if  $\langle \mathbf{u}_i, \mathbf{u}_j \rangle \rightarrow 0 < \pi/2$ .
  - need  $\mathbf{u}_j$ s far from the other as possible; so evenly distributed.
  - $\mathbb{S}^q$ : for larger  $q$ , can get larger angles between  $\mathbf{u}_j$ 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 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  $p$ , only approximate solution,  $j$ th anchor point:

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

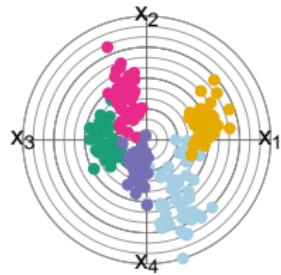
$$u_{j2} = \sin(2\pi j \varphi^{-1}) \sqrt{1 - u_{j3}^2},$$

$$u_{j3} = \frac{2j-1}{p} - 1,$$

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

- distributes anchor points along generative spiral on  $S^2$ , with consecutive points as separated as possible, satisfies "well-separation" property (Saff & Kuijlaars '97).

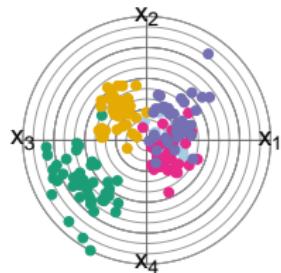
# Illustrative 4D Examples



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

Viz3D

RadViz3D



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

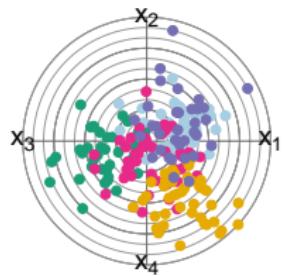
# Higher-dimensional Datasets

- Displaying  $p$  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
- Will 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.
- MRP to maximize separation between groups (in projected space) relative to total variability.

# Maximum-Ratio Projection (MRP)

- **Objective:** Obtain orthogonal projections that maximize between-groups total SSCP relative to total SSCP.
  - Let  $\mathbf{T}$ ,  $\mathbf{W}$  be (p.d.) total & between-groups corrected SSCP.
    - $\hat{\mathbf{v}}_j = \mathbf{T}^{-\frac{1}{2}} \hat{\mathbf{w}}_j / \| \mathbf{T}^{-\frac{1}{2}} \hat{\mathbf{w}}_j \|$ ,  $j = 1, 2, \dots, k$  where  $\hat{\mathbf{w}}_j, j = 1, 2, \dots, k$  are, in decreasing order, the  $k$  largest eigenvalues of  $\mathbf{T}^{-1/2} \mathbf{B} \mathbf{T}^{-1/2}$ .
    - $k \leq G - 1$ , chosen by scree plot/quality of display
    - $G \leq 4$  needs  $4 - G + 1$  more projections w/ null contribution
    - needs p.d.  $\mathbf{T}$ , does not hold if  $p > \min n_g$
- **Solution:** Obtain PCs (orthogonal  $\mathbf{V}_g$ ) for each group
  - Find orthogonal  $\mathbf{W}$  closest to all  $\mathbf{V}_g$ 
    - Project  $\mathbf{X}$  with  $\mathbf{W}$  and then obtain MRP
- Display with RadViz3D

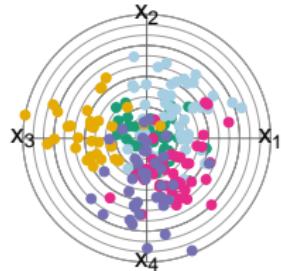
# Illustrative 500D Examples



RadViz2D

Viz3D

RadViz3D

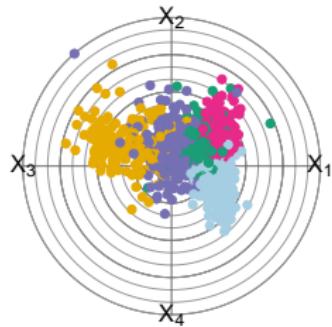


RadViz2D

# 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 Gaussianized Distributional Transform (GDT)
    - marginal application of GDT specifies distribution on  $\mathbf{X}$  with desired marginal and correlation structure.
- GDT standardizing transform, more stringent than usual affine 0-mean, unit-variance inducing transform
  - GDT matches all marginal quantiles to  $N(0, 1)$
  - Apply to skewed datasets or with unclear marginals
- Can use MRP and visualize with RadViz3D

# Applications: Gamma Ray Bursts Dataset



RadViz2D

Viz3D  
Groups ● 1 ● 2 ● 3 ● 4 ● 5

RadViz3D

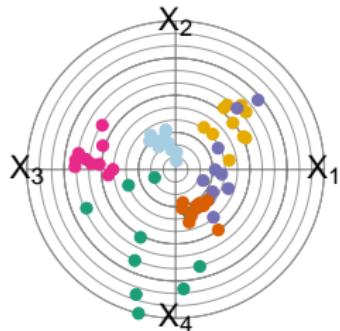
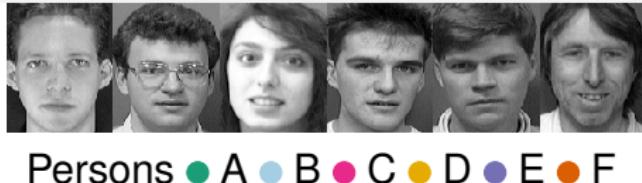
- Heavily skewed attributes, so GDT appropriate
- Results indicate 5 overlapping clusters
  - some suggestion of 2, 3 super-types of GRBs

## Applications: Six Faces



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

# Applications: Six Faces Dataset



RadViz2D

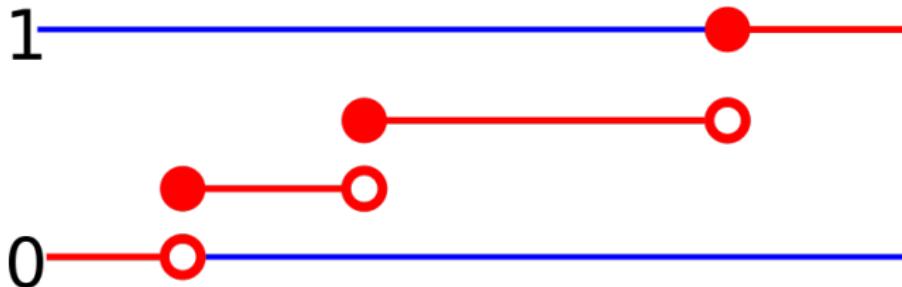
Viz3D

RadViz3D

- use of transforms makes marginal unclear; use GDT
- RadViz2D identifies B, C; Viz3D also A, B

## Datasets with Discrete Attributes

- For discrete-valued RV  $X$ , CDF  $F_X(X) \not\sim U(0, 1)$  because of discreteness.
  - GDT not applicable (yet!)



- Note that the CDF is only right continuous
- Solution proposed by Rüschedendorf (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

$$U = F(X, V) \text{ where}$$

$F(x, \lambda) \doteq P(X < x) + \lambda P(X = x) = F_X(x-) + \lambda[F_X(x) - F_X(x-)]$  is the generalized CDF of  $X$ .

## Theorem

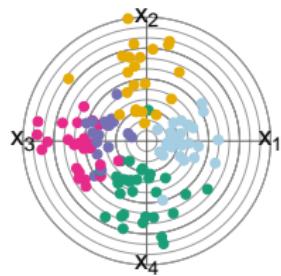
Let  $U = F(X, V)$  be the generalized distributional transform of  $X$ . Then

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

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
  - apply MRP and visualize via RadViz3D

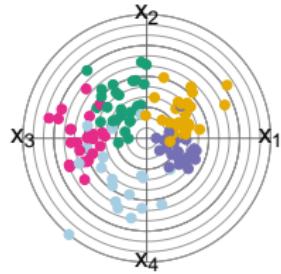
# Illustration: Simulated Binary Datasets



RadViz2D

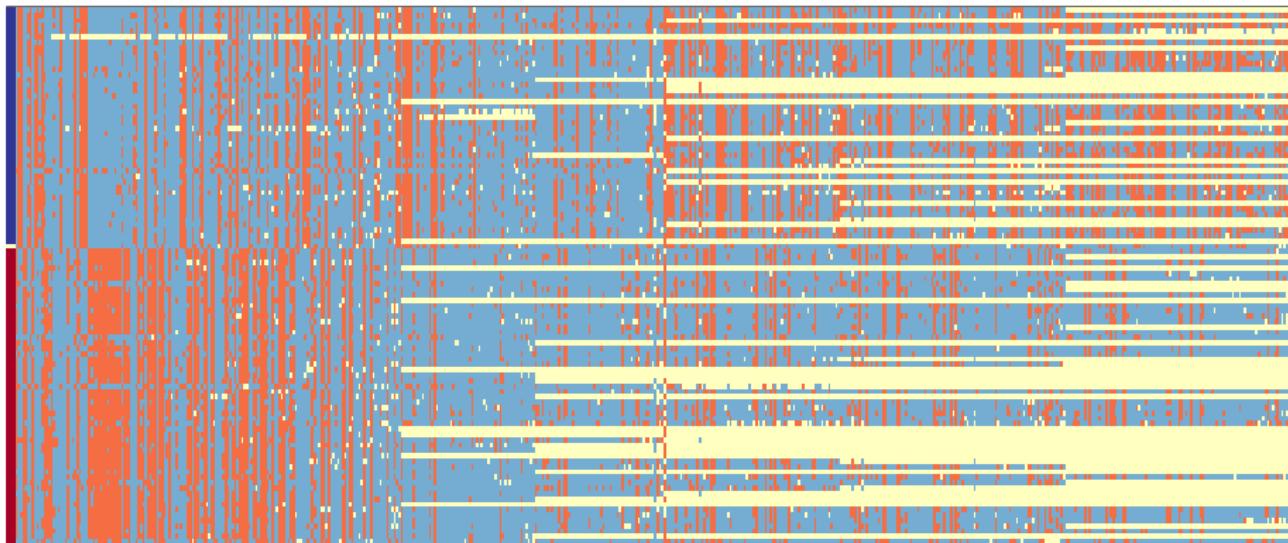
Viz3D

RadViz3D



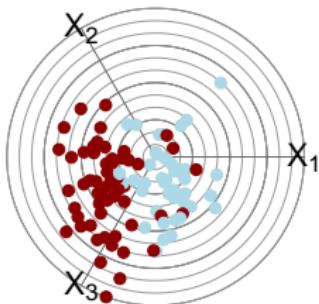
RadViz2D

# 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



RadViz3D

- RadViz3D separates out the two parties best
- $G = 2$  so only 1 MRP with +ve eigenvalue
  - RadViz2D/Viz3D, RadViz3D needs 2, 3 more projections
    - 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)

# Applications: Handwritten Indic Scripts

ଆମ୍ବିନ୍ଦୁ କାହିଁଥାର ଅମର ଜାଗିତ୍ର ହେଲା । ଯାହା ତାଙ୍କୁ,

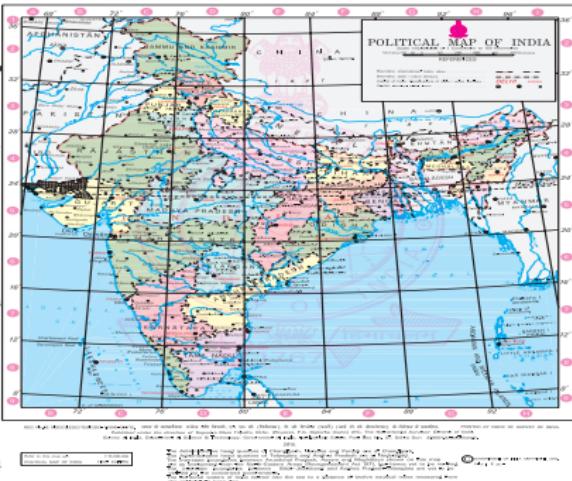
ଏହା କାହିଁଥାର ଯାଏବି ପିଲାରୀ ଦି ଖାଲୀ ଦିଲେ ଏହା ପରା ଧାରିଲେବା-

ଶବ୍ଦରେ ଏହା କାହିଁଥାର କା ମରା ହେବା । ଅଗଜ୍ ଜୋ

ଦୁର୍ଗତ୍ତ ଯେବିଦି କେଲନେବଳ ଅରଧ ଦୁଇଟିରେ ଯନ୍ମ

ଅନ୍ତର୍ଭାବରେ କାହାରେବେଳେ କୁବିତ୍ର କାହାରେବେଳେ

ଅନ୍ତର୍ଭାବରେ କାହାରେବେଳେ - (ଅନ୍ତର୍ଭାବରେ କାହାରେବେଳେ)

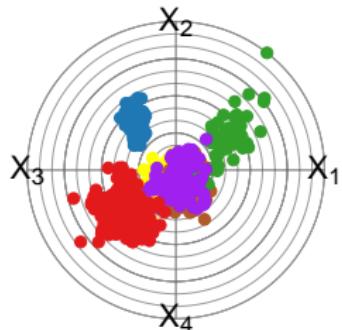


- Handwritten scripts from Bangla, Gujarati, Gurmukhi, Kannada, Malayalam, Urdu

- 116 mixed features (Obaidullah et al, 2017)

(Map Acknowledgment: Surveyor-General of India)

# Applications: Handwritten Indic Scripts



RadViz2D

Viz3D

RadViz3D

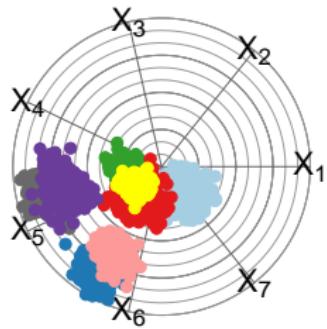
● Bangla ● Gujarati ● Gurmukhi ● Kannada ● Malayalam ● Urdu

- Viz3D (lesser extent RadViz2D) separates Urdu, Kannada and Gujarati, not the other 3 languages
- RadViz3D best in clarifying 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 # samples) 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



- RadViz3D indicates clear separation between organs
  - prostate and stomach have some marginal overlap.
- RadViz2D, Viz3D poorer at separating organs

# 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 available <https://github.com/fanne-stat/radviz3d>
  - manuscript eventually available at arXiv
- 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