

Multivariate Rank Based Methods For Classification

Biman Chakraborty

School of Mathematics, University of Birmingham,
Edgbaston, Birmingham B15 2TT, United Kingdom

Multi-Disciplinary Research in Data Mining - BTG Away Day
20-21 September, 2010

Univariate Sign Function

$$\text{sign}(x) = \begin{cases} -1 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \\ +1 & \text{if } x > 0 \end{cases}$$

$$= \begin{cases} \frac{x}{|x|} & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

$$= \begin{cases} \frac{d}{dx}|x| & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

Univariate Ranks

- Centered Rank:

$$\text{Rank}(x) = \frac{1}{n} \sum_{i=1}^n \text{sign}(x - X_i).$$

- Note that: $E[\text{Rank}(x)] = 2F(x) - 1$.
- $-1 \leq \text{Rank}(x) \leq 1$
- $\text{Rank}(x) = 0 \implies x$ is sample median.

Multivariate Signs:

Consider $\mathbf{x} \in \mathbb{R}^d$

$$\begin{aligned} \text{Sign}(\mathbf{x}) &= \begin{cases} \frac{\partial}{\partial \mathbf{x}} \|\mathbf{x}\| & \text{if } \mathbf{x} \neq \mathbf{0} \\ \mathbf{0} & \text{if } \mathbf{x} = \mathbf{0} \end{cases} \\ &= \begin{cases} \frac{\mathbf{x}}{\|\mathbf{x}\|} & \text{if } \mathbf{x} \neq \mathbf{0} \\ \mathbf{0} & \text{if } \mathbf{x} = \mathbf{0} \end{cases} \end{aligned}$$

where $\|\mathbf{x}\| = \{|x_1|^2 + \dots + |x_d|^2\}^{1/2}$.

Multivariate Ranks:

- Data: $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \in \mathbb{R}^d$.

-

$$\text{Rank}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \text{Sign}(\mathbf{x} - \mathbf{X}_i).$$

- $\|\text{Rank}(\mathbf{x})\| \leq 1$
- $\text{Rank}(\mathbf{x}) = \mathbf{0} \implies \mathbf{x}$ is the spatial median.
- Population Version:

$$R_F(\mathbf{x}) = E_F[\text{Sign}(\mathbf{x} - \mathbf{X})], \quad \mathbf{X} \sim F.$$

Relation With Quantiles:

- Univariate Case:
 - $Rank(x) = u$, $-1 \leq u \leq 1$ then x is the $(u + 1)/2$ -th quantile.
- Multivariate Case: (Chakraborty, 2001)
 - $Rank(\mathbf{x}) = \mathbf{u}$, $\|\mathbf{u}\| < 1$ then \mathbf{x} is the \mathbf{u} -th geometric quantile.
- A Related Depth Function:
 - $D(\mathbf{x}) = 1 - \|Rank(\mathbf{x})\|$.
 - Depth of the median = 1.

A Simple Classification Rule

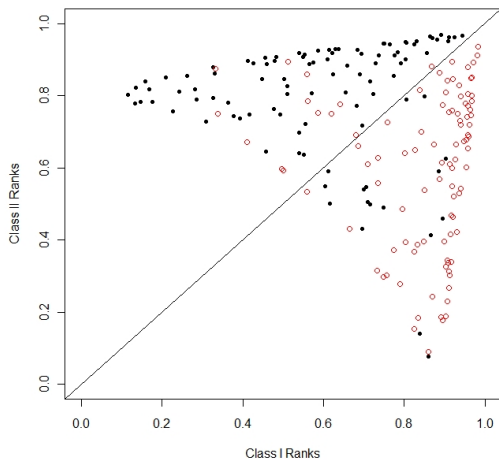
- Two-class Problem:

- Let C_1 and C_2 be the two classes of population with the d -dimensional feature vector \mathbf{X} having distributions F and G , respectively.
- **Rule:** Classify $\mathbf{x} \in \mathbb{R}^d$ to C_1 if $\|R_F(\mathbf{x})\| < \|R_G(\mathbf{x})\|$ and classify \mathbf{x} to C_2 otherwise.

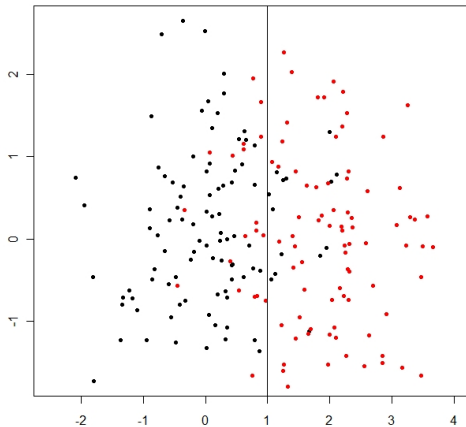
- K -class Problem:

- Let the classes C_1, \dots, C_K have distributions F_1, \dots, F_K .
- **Rule:** Classify \mathbf{x} to C_i if $\|R_{F_i}(\mathbf{x})\| = \min_{1 \leq j \leq K} \|R_{F_j}(\mathbf{x})\|$.

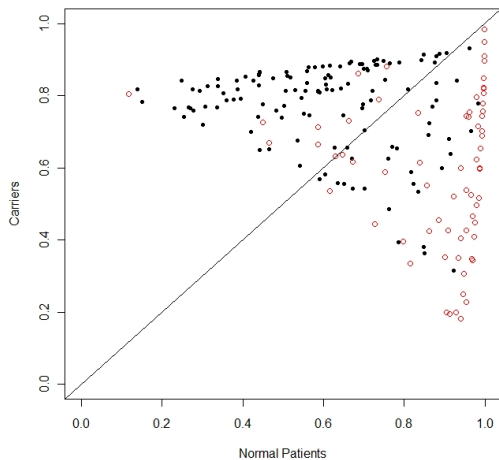
A Simulated Example: Rank-Rank Plot



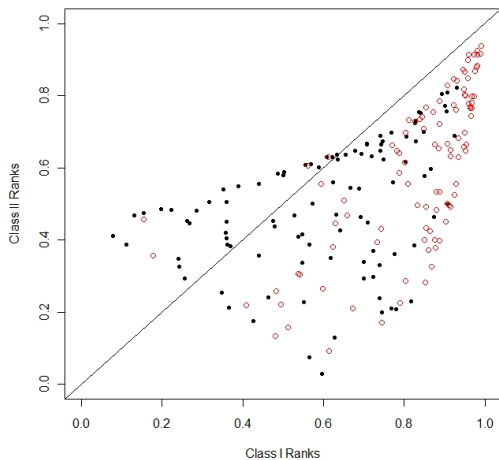
The Example: In The Feature Space



Blood Measurements Data



A Difference in Scale Example



Comments

- The previous rule is Bayes optimal for elliptically symmetric distributions with a location shift.
- It does not perform well for difference in scales.
- We may consider rules like: Classify to C_1 if $\|R_F(\mathbf{x})\| < a\|R_G(\mathbf{x})\| + c$ for some optimal choice of a and c .
- In general, one can consider: Classify to C_1 if $\|R_F(\mathbf{x})\| < h(\|R_G(\mathbf{x})\|)$ for some function h .

Comments

- Similar ideas can be implemented for clustering.
- We can also consider some local version of multivariate ranks.
- Similar ideas can be explored with other notions of data depths.
- **Problem:** Really large and high dimensional dataset.
- **Problem:** Feature vector containing discrete variables.