



# On Some Fast And Robust Classifiers For HDLSS Data

Sarbojit Roy<sup>1</sup> Jyotishka Ray Choudhury<sup>2</sup> Subhajit Dutta<sup>1</sup>

<sup>1</sup>Indian Institute of Technology Kanpur <sup>2</sup>Indian Statistical Institute Kolkata



## Classification in High Dimension, Low Sample Size Settings

Suppose  $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$  and  $\mathbf{Y}_j = (Y_{j1}, \dots, Y_{jp})^\top$  are i.i.d. random vectors from distribution functions  $\mathbf{F}_1$  and  $\mathbf{F}_2$ , respectively, for  $1 \leq i \leq n_1$  and  $1 \leq j \leq n_2$ . Let  $X = X_1 \cup X_2$  denote the training sample of size  $n = n_1 + n_2$ , where  $X_1 = \{\mathbf{X}_1, \dots, \mathbf{X}_{n_1}\}$  and  $X_2 = \{\mathbf{Y}_1, \dots, \mathbf{Y}_{n_2}\}$ . The class prior probabilities  $0 < \pi_1, \pi_2 < 1$  satisfy  $\pi_1 + \pi_2 = 1$ .

Given the training sample  $X$ , our aim is to develop a classifier  $\delta$  such that its misclassification probability  $\Delta$  goes to zero under fairly general conditions in the high dimension, low sample size (HDLSS) regime, where  $n$  is held fixed while  $p \rightarrow \infty$ .

## Limitations of Existing Classifiers

- $\mu_1 = E[\mathbf{X}]$ ,  $\mu_2 = E[\mathbf{Y}]$ ,  $\Sigma_1 = \text{Cov}[\mathbf{X}]$  and  $\Sigma_2 = \text{Cov}[\mathbf{Y}]$ .
- Define  $\nu^2 = \lim_{p \rightarrow \infty} \frac{1}{p} \|\mu_1 - \mu_2\|^2$  and  $\sigma_j^2 = \lim_{p \rightarrow \infty} \frac{1}{p} \text{trace}(\Sigma_j)$  for  $j = 1, 2$ .
- Existing classifiers yield *perfect classification* (i.e.,  $\Delta \rightarrow 0$  as  $p \rightarrow \infty$ ) if
  - $\nu^2 > |\sigma_1^2 - \sigma_2^2|$  for the nearest neighbor (NN) classifier, average distance (AVG) classifier, support vector machines (SVM) [2].
  - $\nu^2 > 0$ , or  $\sigma_1^2 \neq \sigma_2^2$  for the scale adjusted AVG (SAVG) classifier [1].
- In HDLSS settings, behavior of the existing classifiers is governed by the constants  $\nu^2$ ,  $\sigma_1^2$  and  $\sigma_2^2$ .

## Our Contribution

We propose classifiers that are **robust**, **computationally fast**, **free from tuning parameters** and have **strong theoretical properties**.

## A Robust and Tuning-free Classifier

Define  $h(u, v) = \sin^{-1} \left( (1 + uv) / \sqrt{(1 + u^2)(1 + v^2)} \right) / 2\pi$  for  $u, v \in \mathbb{R}$  and

$$\bar{h}(\mathbf{u}, \mathbf{v}) = \frac{1}{p} \sum_{k=1}^p h(u_k, v_k) \text{ for } \mathbf{u} = (u_1, \dots, u_p)^\top, \mathbf{v} = (v_1, \dots, v_p)^\top \in \mathbb{R}^p.$$

$$\bar{T}_{11} = \sum_{i < j} \frac{\bar{h}(\mathbf{X}_i, \mathbf{X}_j)}{n_1(n_1 - 1)}, \bar{T}_{12} = \sum_{i,j} \frac{\bar{h}(\mathbf{X}_i, \mathbf{Y}_j)}{n_1 n_2}, \bar{T}_{22} = \sum_{i < j} \frac{\bar{h}(\mathbf{Y}_i, \mathbf{Y}_j)}{n_2(n_2 - 1)},$$

$$\bar{T}_1(\mathbf{z}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \bar{h}(\mathbf{X}_i, \mathbf{z}), \bar{T}_2(\mathbf{z}) = \frac{1}{n_2} \sum_{i=1}^{n_2} \bar{h}(\mathbf{Y}_i, \mathbf{z}), \bar{L}_j(\mathbf{z}) = \bar{T}_{jj} - 2\bar{T}_j(\mathbf{z}) \text{ for } j = 1, 2.$$

- Discriminant:**  $\bar{L}(\mathbf{z}) = \bar{L}_2(\mathbf{z}) - \bar{L}_1(\mathbf{z})$ .
- Classifier:**  $\delta_1(\mathbf{z}) = \begin{cases} 1, & \text{if } \bar{L}(\mathbf{z}) > 0, \\ 2, & \text{otherwise.} \end{cases}$

## A measure of distance between $\mathbf{F}_1$ and $\mathbf{F}_2$

- $h$  is a bounded function and free of parameters.
- Define  $\bar{\tau}_p = E[\bar{h}(\mathbf{X}_1, \mathbf{X}_2)] + E[\bar{h}(\mathbf{Y}_1, \mathbf{Y}_2)] - 2E[\bar{h}(\mathbf{X}_1, \mathbf{Y}_1)]$  for  $p \geq 1$ .
- Clearly,  $\bar{\tau}_p \geq 0$ . Equality holds iff  $F_{1k} = F_{2k}$  for all  $1 \leq k \leq p$  where  $F_{1k}$  and  $F_{2k}$  are one-dimensional marginals of  $\mathbf{F}_1$  and  $\mathbf{F}_2$ , respectively (see [3]).

- $\bar{\tau}_p$  is a measure of distance between  $\mathbf{F}_1$  and  $\mathbf{F}_2$ .
- Now,  $E[\bar{L}(\mathbf{Z}) \mid \mathbf{Z} \sim \mathbf{F}_1] = \bar{\tau}_p \geq 0$ , while  $E[\bar{L}(\mathbf{Z}) \mid \mathbf{Z} \sim \mathbf{F}_2] = -\bar{\tau}_p \leq 0$ .

## Limitations of $\bar{\tau}_p$

Define  $\bar{\tau}_p(1, 1) = E[\bar{h}(\mathbf{X}_1, \mathbf{X}_2)]$ ,  $\bar{\tau}_p(2, 2) = E[\bar{h}(\mathbf{Y}_1, \mathbf{Y}_2)]$  and  $\bar{\tau}_p(1, 2) = E[\bar{h}(\mathbf{X}_1, \mathbf{Y}_1)]$ .

Observe that  $\bar{\tau}_p = \{\bar{\tau}_p(1, 1) - \bar{\tau}_p(1, 2)\} + \{\bar{\tau}_p(2, 2) - \bar{\tau}_p(1, 2)\}$ .

- Suppose  $\bar{\tau}_p(1, 1) < \bar{\tau}_p(1, 2) < \bar{\tau}_p(2, 2)$ . Then, the value of  $\bar{\tau}_p$  may become small.
- An improved dissimilarity index:  $\bar{\psi}_p = \{\bar{\tau}_p(1, 1) - \bar{\tau}_p(1, 2)\}^2 + \{\bar{\tau}_p(2, 2) - \bar{\tau}_p(1, 2)\}^2$ .
  - Squaring the terms before addition eliminates the possibility of cancellations.
  - Further,  $\bar{\psi}_p = 0$  iff  $F_k = G_k$  for all  $1 \leq k \leq p$ .

## A Classifier Based on $\bar{\psi}_p$

Define  $\bar{T} = \bar{T}_{11} - 2\bar{T}_{12} + \bar{T}_{22}$ .

- Discriminant:**  $\bar{\theta}(\mathbf{z}) = \bar{T}\bar{L}(\mathbf{z})/2 + \{\bar{T}_{22} - \bar{T}_{11}\} \{\bar{L}_1(\mathbf{z}) + \bar{L}_2(\mathbf{z}) + 2\bar{T}_{12}\}/2$ .
- Classifier:**  $\delta_2(\mathbf{z}) = \begin{cases} 1, & \text{if } \bar{\theta}(\mathbf{z}) > 0, \\ 2, & \text{otherwise.} \end{cases}$

## Asymptotic Behavior in HDLSS Settings

Suppose  $\mathbf{U} = (U_1, \dots, U_p)^\top$  and  $\mathbf{V} = (V_1, \dots, V_p)^\top$  are two independent vectors such that  $\mathbf{U} \sim \mathbf{F}_j$  and  $\mathbf{V} \sim \mathbf{F}_{j'}$  for  $j, j' \in \{1, 2\}$ . Let us assume the following:

- Weak dependence among the component variables:

$$(A1) \sum_{1 \leq k < k' \leq p} \text{Corr}(h(U_k, V_k), h(U_{k'}, V_{k'})) = o(p^2).$$

- Assumption (A1) is trivially satisfied if the component variables of the underlying distributions are independently distributed.
- It continues to hold when the components have weak dependence among them. For example, (A1) is satisfied when  $\{h(U_k, V_k), k \geq 1\}$  has the  $\rho$ -mixing property.
- If assumption (A1) is satisfied, then we have the following:

$\mathbf{Z} \sim \mathbf{F}_1$	$ \bar{L}(\mathbf{Z}) - \bar{\tau}_p  \xrightarrow{P} 0$ and $ \bar{\theta}(\mathbf{Z}) - \bar{\psi}_p  \xrightarrow{P} 0$ as $p \rightarrow \infty$
$\mathbf{Z} \sim \mathbf{F}_2$	$ \bar{L}(\mathbf{Z}) + \bar{\tau}_p  \xrightarrow{P} 0$ and $ \bar{\theta}(\mathbf{Z}) + \bar{\psi}_p  \xrightarrow{P} 0$ as $p \rightarrow \infty$ .

- Asymptotic separability of  $\mathbf{F}_1$  and  $\mathbf{F}_2$ :

$$(A2) \liminf_{p \rightarrow \infty} \bar{\tau}_p > 0.$$

- If the component variables are identically distributed, then (A2) is satisfied.
- (A2) also implies that  $\liminf_{p \rightarrow \infty} \bar{\psi}_p > 0$ .

## Theorem 1: Perfect Classification

If (A1) and (A2) are satisfied, then for any  $\pi_1 > 0$ , we have  $\Delta_1 \rightarrow 0$  and  $\Delta_2 \rightarrow 0$  as  $p \rightarrow \infty$ .

## Relative Performance of $\delta_1$ and $\delta_2$

- Both  $\delta_1$  and  $\delta_2$  yield *perfect classification* under the same set of assumptions.
- We now provide a set of sufficient conditions under which one classifier outperforms the other. First, let us assume the following:
 

(A3) There exists a  $p_0 \in \mathbb{N}$  such that  $\bar{\tau}_p(1, 2) > \min\{\bar{\tau}_p(1, 1), \bar{\tau}_p(2, 2)\}$  for all  $p \geq p_0$ .
- If assumption (A3) is satisfied, then either of  $\bar{\tau}_p(1, 1) - \bar{\tau}_p(1, 2)$  and  $\bar{\tau}_p(2, 2) - \bar{\tau}_p(1, 2)$  is positive, while the other quantity is negative.

## Theorem 2: Ordering Between $\Delta_1$ and $\Delta_2$

If assumptions (A1) – (A3) are satisfied, then there exists an integer  $p'_0$  such that  $\Delta_2 \leq \Delta_1$  for all  $p \geq p'_0$ .

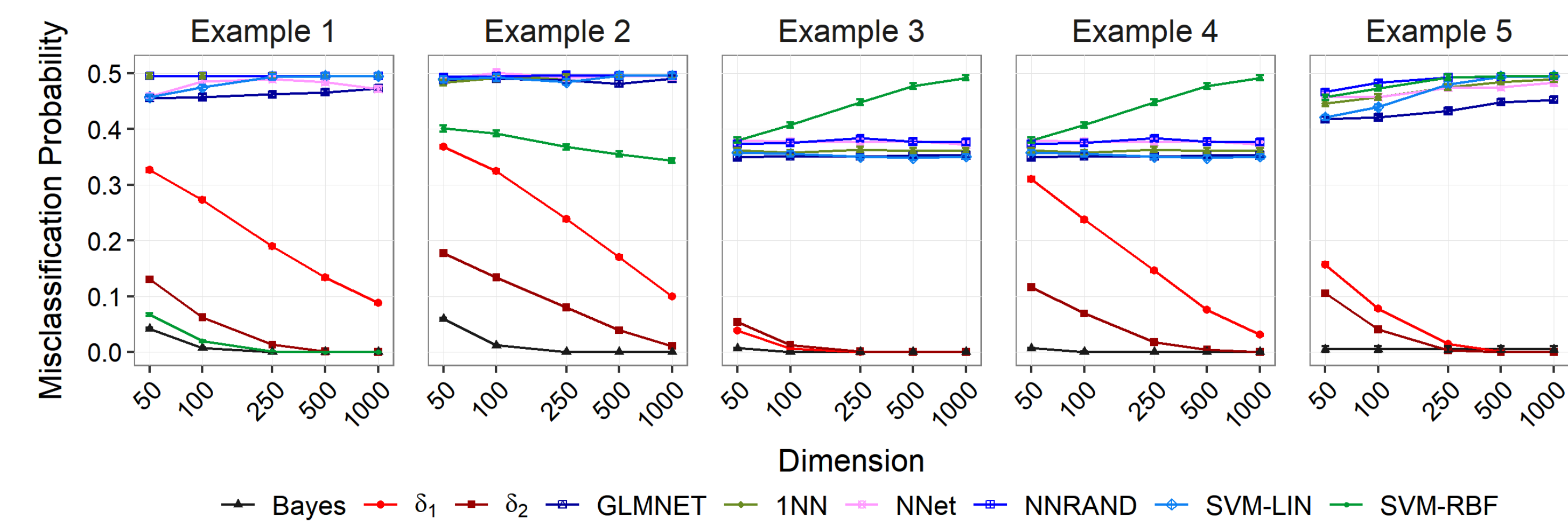
If the inequality in (A3) is inverted, then the ordering of  $\Delta_1$  and  $\Delta_2$  is reversed.

## Simulation Study

Fix  $1 \leq k \leq p$ . Now, consider the following examples:

Example	$(\nu, \sigma_1^2, \sigma_2^2)$	$\bar{T}_{12} > \min\{\bar{T}_{11}, \bar{T}_{22}\}$
1. $X_{1k} \stackrel{i.i.d.}{\sim} N(1, 1)$ and $Y_{1k} \stackrel{i.i.d.}{\sim} N(1, 2)$	$\nu^2 <  \sigma_1^2 - \sigma_2^2 $	True
2. $X_{1k} \stackrel{i.i.d.}{\sim} N(0, 3)$ and $Y_{1k} \stackrel{i.i.d.}{\sim} t_3$	$\nu^2 = \sigma_1^2 - \sigma_2^2 = 0$	True
3. $X_{1k} \stackrel{i.i.d.}{\sim} C(0, 1)$ and $Y_{1k} \stackrel{i.i.d.}{\sim} C(1, 1)$	do not exist	False
4. $X_{1k} \stackrel{i.i.d.}{\sim} C(0, 1)$ and $Y_{1k} \stackrel{i.i.d.}{\sim} C(0, 2)$	do not exist	True
5. $X_{1k} \stackrel{i.i.d.}{\sim} \text{Par}(1, 1)$ and $Y_{1k} \stackrel{i.i.d.}{\sim} \text{Par}(1.25, 1)$	do not exist	True

$N(\mu, \sigma)$ : Gaussian distribution with mean  $\mu \in \mathbb{R}$  and variance  $\sigma > 0$ .  
 $t_\alpha$ : the Student's  $t$ -distribution with  $\alpha > 0$  degrees of freedom.  
 $C(\mu, \sigma)$ : Cauchy distribution with location  $\mu \in \mathbb{R}$  and scale  $\sigma > 0$ .  
 $\text{Par}(\theta, s)$ : Pareto distribution with  $\theta > 0$  and scale  $s > 0$ .



## References

- [1] Yao-Ban Chan and Peter Hall. Scale adjustments for classifiers in high-dimensional, low sample size settings. *Biometrika*, 96(2):469–478, 2009.
- [2] Peter Hall, J. S. Marron, and Amnon Neeman. Geometric representation of high dimension, low sample size data. *Journal of the Royal Statistical Society Series B*, 67(3):427–444, 2005.
- [3] Zhimei Li and Yaowu Zhang. On a projective ensemble approach to two sample test for equality of distributions. In *International Conference on Machine Learning*, pages 6020–6027. PMLR, 2020.