

# HS-, HiS-, HeRS-Boot: Heterogeneity-Stratified Bootstrap, Enhancing Substance Recognition for Assistive technologies in Real-World Applications

Pertami J. Kunz and Syrine ben Abid  
pertami.kunz@ieee.org  
syrinebenabid@ieee.org  
Technische Universität Darmstadt  
Germany

## ABSTRACT

We introduce the Heterogeneity-Stratified, the improved Heterogeneity-Stratified, and the Heterogeneity Ratio-Stratified Bootstrap oversampling methods (HS-, HiS, and HeRS-Boot, respectively), which assign higher resampling probabilities to sample points in less homogeneous regions. We demonstrate its advantage in the case of training a detector by oversampling the under-represented class in an imbalanced data set. We took a case study of spoiled food and allergen detectors in form of an electronic nose. Results demonstrate the effectiveness and generalizability of our method across different sensors, highlighting its potential for real-world applications and positive impact on daily life.

## CCS CONCEPTS

• Mathematics of computing → Bootstrapping.

## KEYWORDS

Oversampling, Bootstrap, Out of Bag Bootstrap, Detection, Classification, IoT

## ACM Reference Format:

Pertami J. Kunz and Syrine ben Abid. 2023. HS-, HiS-, HeRS-Boot: Heterogeneity-Stratified Bootstrap, Enhancing Substance Recognition for Assistive technologies in Real-World Applications. In *Proceedings of WomENcourage ACM 2023*. ACM, New York, NY, USA, 2 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Anosmic or visually impaired people may fail to recognise spoiled foods or allergend in their fridge or pantry. A smart detector that can be trained to recognise such hazards is needed. Intelligently choosing a method to train the algorithm may reduce dependencies on the complexity and the amount of data collected. It is often the case that the class distribution in a dataset is not equal. A balanced training dataset may be created with an artificially equal class distribution by oversampling the instances in the minority class. In

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*WomENcourage ACM 2023, Sept 20–22, 2023, Trondheim, Norway*

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

this paper, we introduce three oversampling methods based on the stratified resampling with replacement, or stratified bootstrapping.

## 2 METHODOLOGY

For the dataset, the gas compositions of the following specimens were measured The BME688 (Bosch Sensortec) sensor: (1) Fresh Chicken: a piece of fresh, raw chicken, (2) Yoghurt: fresh yoghurt (3) Beef: a piece of fresh raw beef, (4) Coffee: a handful of coffee beans (5) Mix : a piece of spoiled raw chicken, mixed with some fresh vegetables (6) Rotten Chicken : a piece of spoiled raw chicken. The measurement setup and the plot of the measurements are shown in Figure 1. The proposed algorithms are summarised in Algorithm 1.

### Algorithm 1 Oversampling with the HS-, HiS, or HeRS-Boot [1, 2]

- Step 1** Divide the sample domain into  $K$  grids.
- Step 2** For each grid  $k$ ,  $k = 1, \dots, K$ , calculate the heterogeneity measure  $H(C, k)$  and modify the resampling probability distribution  $p_c(k)$  for each class  $c$  in cluster  $k$  (Table 1).
- Step 3** Oversampling: Sample  $N_0$  instances with replacement from Class 0 (non-target) and likewise  $N_0$  instances from Class 1 (target) with the updated distribution  $p'_c$ .

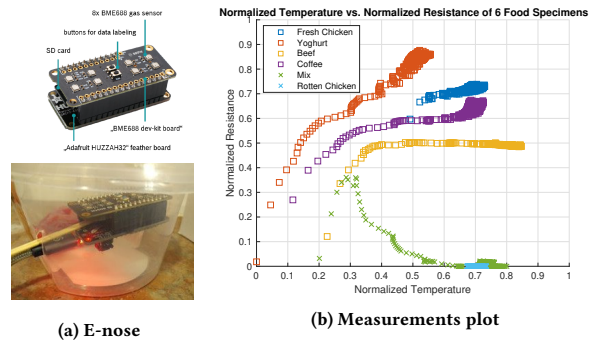


Figure 1: Training the electronic nose with different specimens.

## 3 RESULTS AND DISCUSSION

Tables 2 and 3 show the detection metrics averaged over 500 trials, where TP is the true positives, TN the true negatives, FP the false positive, and FN the false negatives: (1) **Sensitivity (SN)** or true positive rate (TPR) or recall (REC) =  $TP / (TP + FN)$ , (2) **Specificity (SP)** or true negative rate (TNR) =  $TN / (TN + FP)$ , (3) **Accuracy (ACC)** =  $(TP + TN) / (TP + TN + FN + FP)$ , (4) **Error rate (ERR)** =

**Table 1: Summary of Algorithms, where  $\gamma$  is a hyperparameter constant that determines the influence of the heterogeneity,  $a_{c,k}$  is the number of samples in grid  $k$  that belongs to class  $c$ , and the denominator is to make sure that the sum equals 1.**

HSBoot, HiSBoot		HeRSBoot	
$H(C, k) = -\sum_{c=1}^C \frac{a_{c,k}}{N} \log\left(\frac{a_{c,k}}{\sum_1^C a_{c,k}}\right)$		$H(C, k) = -\sum_{c=1}^C \frac{p_{c,k}}{K} \log(p_{c,k})$	
HSBoot		HiSBoot, HeRSBoot	
$p'_c(k) = \frac{\frac{1}{N_c} + \gamma H(C, k)}{1 + \gamma \sum_k a_{c,k} H(C, k)}$	$p'_c(k) = \begin{cases} \frac{\frac{1}{N_c} + \gamma H(C, k)}{1 + \gamma \sum_k a_{c,k} H(C, k)} & \text{for target class} \\ p_c(k) & \text{for non target class} \end{cases}$		

(FP + FN) / (TP + TN + FN + FP), (5) **Precision (PREC)** = TP / (TP + FP), and (6) **F1 score (F1)** = 2 \* PREC \* REC / (PREC + REC).

The Naive Bayes, Quadratic Discriminant, and Neural Network algorithms were used as classifiers, where in each trial, the hyperparameters were decided with leave 0.3 out cross validation. We can see from the tables, the proposed methods with different  $\lambda$  improve the traditional training set selection method.

**Table 2: Spoiled Food Detection Results**

(a) Naive Bayes

	Random	HiSBoot <sup>2</sup> ,			HeRSBoot <sup>2</sup>		
		$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$
SN	.9899	.9885	.9922	.9902	.9911	<b>.9924</b>	<b>.9924</b>
SP	.9924	<b>.9953</b>	.9945	.9927	.9946	.9923	.9917
ACC	.9916	.9930	<b>.9937</b>	.9919	.9934	.9923	.9920
ERR	.0084	.0070	<b>.0063</b>	.0081	.0066	.0077	.0080
PREC	.9850	<b>.9906</b>	.9891	.9856	.9892	.9848	.9837
F1	.9872	.9894	<b>.9906</b>	.9877	.9901	.9885	.9879

(b) Quadratic Discriminant

	Random	HiSBoot <sup>2</sup> ,			HeRSBoot <sup>2</sup>		
		$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$
SN	.9941	.9916	.9931	<b>.9953</b>	.9925	.9944	.9948
SP	.9913	<b>.9947</b>	.9930	.9901	.9938	.9911	.9917
ACC	.9922	<b>.9937</b>	.9930	.9918	.9933	.9922	.9928
ERR	.0078	<b>.0063</b>	.0070	.0082	.0067	.0078	.0072
PREC	.9831	<b>.9898</b>	.9863	.9807	.9879	.9826	.9837
F1	.9885	<b>.9906</b>	.9896	.9879	.9901	.9884	.9892

(c) Neural Network with 2 hidden layers, with 8 and 4 fully connected outputs for each hidden layer, respectively.

	Random	HiSBoot <sup>2</sup> ,			HeRSBoot <sup>2</sup>		
		$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$
SN	<b>.9869</b>	.9303	.9446	.9309	.9208	.9531	.9352
SP	.9253	.9654	<b>.9664</b>	.9569	.9605	.9535	.9518
ACC	.9459	.9537	<b>.9591</b>	.9482	.9473	.9534	.9462
ERR	.0541	.0463	<b>.0409</b>	.0518	.0527	.0466	.0538
PREC	.9094	.9406	<b>.9429</b>	.9274	.9325	.9244	.9194
F1	.9449	.9652	<b>.9667</b>	.9579	.9607	.9562	.9536

**Table 3: Allergen Detection Results**

(a) Naive Bayes

	Random	HiSBoot <sup>2</sup> ,			HeRSBoot <sup>2</sup>		
		$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$
SN	.6053	.7992	.7583	.7128	.8122	<b>.8157</b>	.7823
SP	.9810	.9762	.9793	<b>.9816</b>	.9744	.9728	.9735
ACC	.9184	.9467	.9425	.9368	<b>.9474</b>	.9466	.9417
ERR	.0816	.0533	.0575	.0632	<b>.0526</b>	.0534	.0583
PREC	.8655	.8723	.8819	<b>.8877</b>	.8653	.8590	.8576
F1	.7020	.8296	.8117	.7870	<b>.8338</b>	.8325	.8140

(b) Quadratic Discriminant

	Random	HiSBoot <sup>2</sup> ,			HeRSBoot <sup>2</sup>		
		$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$
SN	.8057	.7591	.7931	.7509	.7503	.8310	<b>.8474</b>
SP	.9653	.9701	.9697	<b>.9762</b>	.9704	.9653	.9631
ACC	.9387	.9349	.9402	.9387	.9337	.9429	<b>.9438</b>
ERR	.0613	.0651	.0598	.0613	.0663	.0571	<b>.0562</b>
PREC	.8228	.8353	.8401	<b>.8647</b>	.8353	.8273	.8211
F1	.8141	.7952	.8154	.8030	.7904	.8289	<b>.8334</b>

(c) Neural Network with 2 hidden layers, with 8 and 4 fully connected outputs for each hidden layer, respectively.

	Random	HiSBoot <sup>2</sup> ,			HeRSBoot <sup>2</sup>		
		$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 5$
SN	<b>.7252</b>	.6395	.6399	.6276	.6228	.6506	.6474
SP	.9620	.9805	.9811	<b>.9824</b>	.9801	.9768	.9764
ACC	.9225	.9237	<b>.9242</b>	.9233	.9205	.9225	.9216
ERR	.0775	.0763	<b>.0758</b>	.0767	.0795	.0775	.0784
PREC	.8315	.8681	.8707	<b>.8782</b>	.8627	.8496	.8465
F1	<b>.7737</b>	.7548	.7656	.7690	.7524	.7653	.7702

**SELECTED REFERENCES**

- [1] Pertami J. Kunz and Abdelhak M. Zoubir. 2023. Heterogeneity-Stratified Bootstrap Oversampling for Training a Spoiled Food Detector. (2023). Submitted to DSP 2023, Rhodos, Greece.
- [2] Pertami J. Kunz and Abdelhak M. Zoubir. 2023. The improved Heterogeneity and Heterogeneity Ratio-Stratified Bootstrap Oversampling for Training an Allergen Detector. (2023). Draft.