



# Novelty Discovery with Kernel Minimum Enclosing Balls

Rafet Sifa<sup>1(✉)</sup> and Christian Bauckhage<sup>1,2</sup>

<sup>1</sup> Fraunhofer Center for Machine Learning, Sankt Augustin, Germany  
{rafet.sifa, christian.bauckhage}@iais.fraunhofer.de

<sup>2</sup> University of Bonn, Bonn, Germany

**Abstract.** We introduce the idea of utilizing ensembles of Kernel Minimum Enclosing Balls to detect novel datapoints. To this end, we propose a novelty scoring methodology that is based on combining outcomes of the corresponding characteristic functions of a set of fitted balls. We empirically evaluate our model by presenting experiments on synthetic as well as real world datasets.

## 1 Introduction

The notion of novelty discovery (or detection) [6] can be described as a one-class classification problem (a.k.a *data domain description* [11]) aiming to learn certain characteristics of the analyzed datasets to be able to separate novel datapoints. It finds many applications in numerous scientific and engineering areas such as fraudulent activity detection in financial applications or detecting rare events in medical monitoring [6]. Although reservoir computing based approaches [4] have been proposed to a variety of classification and regression problems, to the best of our knowledge, corresponding methods that are oriented to tackle one-class problems are scarce. The main contribution of this work is about utilizing Minimum Enclosing Balls [1] for novelty discovery. Minimum Enclosing Balls (MEBs) fall into the class of unsupervised representation learning methods that can be used to extract important characteristics about the considered datasets [1]. The main idea behind the MEBs is about determining the smallest ball encapsulating the *entire* dataset in the data- or feature space, which can be found by formulating the problem as an inequality constrained convex minimization problem with a dual allowing for invoking the kernel trick and this dual can be solved using dynamical processes from reservoir computing [1].

Our contribution is based on the decisions of a set of Kernel Minimum Enclosing Balls (KMEBs) by introducing a compound novelty score, which can allow for, for instance, a majority voting based detection as decision based on single balls might be limiting for novelty detection. In addition, our methodology

---

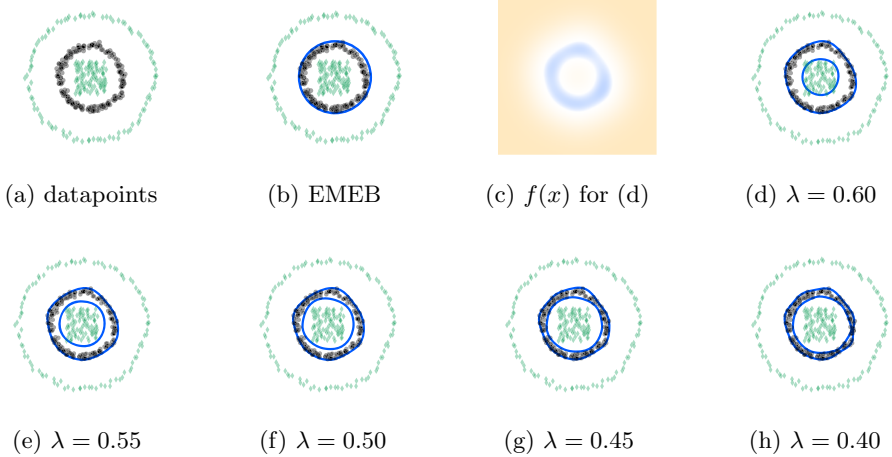
Supported by the Competence Center for Machine Learning Rhine Ruhr (ML2R) which is funded by the Federal Ministry of Education and Research of Germany (grant no. 01—S18038A).

© Springer Nature Switzerland AG 2020

I. S. Kotsireas and P. M. Pardalos (Eds.): LION 14 2020, LNCS 12096, pp. 414–420, 2020.

[https://doi.org/10.1007/978-3-030-53552-0\\_37](https://doi.org/10.1007/978-3-030-53552-0_37)

can be easily implemented in neuromorphic architectures and is capable of dealing with nonlinear patterns due to kernelization [1, 10]. Figure 1 shows an illustrative example explaining our idea about detecting the novel datapoints (green diamond shaped points in Fig. 1a) given a dataset of normal datapoints (black round points in Fig. 1a), which can neither be detected using euclidean Minimum Enclosing Balls (as seen in Fig. 1b) nor considering probabilistic novelty detection such as the deviation from the sample data mean [6]. Instead by considering the characteristic functions of multiple KMEBs (see example in Fig. 1c) with differently scaled Gaussian kernels we can detect all novel points that the considered balls might not individually be capable of capturing (compare the results of Fig. 1d to the others).



**Fig. 1.** A conceptual example illustrating the idea of utilizing Kernel Minimum Enclosing Balls for novelty discovery. (a) shows the data (the inner ring), which is used to compute the ball, and the novel (green diamond) points. It is important to note that neither considering the deviation from the mean vector nor computing the euclidean MEB, which is shown in (b), can in this case isolate the points inside the inner ring. (c) shows a heat-map of the characteristic function from Eq. 7, where colors orange, white and blue respectively indicate positive, zero and negative values. (d–h) shows the dataset and the novel points with the decision boundaries for different Gaussian Kernel scale values  $\lambda$ . Used individually to detect the novel points, the recall values for detecting the novelty are respectively 0.935, 0.995 and 1.000 for the balls in (d), (e) and (f–h). We obtain 1.000 recall when majority-voting over the prediction of the balls (d–h) (i.e. by considering an ensemble of 5 KMEBs with evenly spaced  $\lambda$  values over  $[0.4, 0.6]$ ). (Color figure online)

The remaining of the paper is organized as follows. So as to be self-contained, in Sect. 2 we will formally define the notion of KMEBs, show how we can compute them following a process akin to the ones used in echo state networks and finally show how, once computed, the support vectors of balls can be used to

characterize the interior of the fitted balls. Following that in Sect. 3 we will introduce a new novelty scoring methodology based on the characteristic functions for novelty discovery. In Sect. 4 we will present empirical results to evaluate our approach using real world datasets and in Sect. 5 we will conclude our work.

## 2 An Overview of Kernel Minimum Enclosing Balls

Given a set of  $m$ -dimensional data points  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  (for  $\mathbf{x}_i \in \mathbb{R}^m$ ) that are grouped into a column data matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ , we aim to find the  $m$ -ball  $\mathcal{B}(\mathbf{c}, r)$  containing each of the given data points in  $\mathcal{X}$ , where  $\mathbf{c} \in \mathbb{R}^m$  and  $r \in \mathbb{R}$  are respectively the center and the radius of  $\mathcal{B}$ . Finding MEBs can be cast as a convex optimization problem

$$\begin{aligned} \mathbf{c}_*, r_* &= \underset{\mathbf{c}, r}{\operatorname{argmin}} \quad r^2 \\ \text{s. t.} \quad & \|\mathbf{x}_i - \mathbf{c}\|^2 - r^2 \leq 0 \quad i \in [1, \dots, n]. \end{aligned} \tag{1}$$

Upon evaluating the Lagrangian and the KKT conditions, the negated dual of (1), allows for the kernel trick (as the data *only* occurs in form of inner products [1]) and can be written as the minimization problem

$$\begin{aligned} \boldsymbol{\mu}_* &= \underset{\boldsymbol{\mu}}{\operatorname{argmin}} \quad \boldsymbol{\mu}^\top \mathbf{K} \boldsymbol{\mu} - \boldsymbol{\mu}^\top \mathbf{k} \\ \text{s. t.} \quad & \sum_{i=1}^n \mu_i = 1 \quad \wedge \quad \mu_j \geq 0 \quad \forall j \in [1, \dots, n], \end{aligned} \tag{2}$$

where  $\mathbf{K} \in \mathbb{R}^{n \times n}$  is a kernel matrix,  $\mathbf{k}$  contains its diagonal (i.e.  $\mathbf{k} = \operatorname{diag}[\mathbf{K}]$ ) and  $\boldsymbol{\mu} \in \mathbb{R}^n$  contains Lagrange multipliers. The kernel matrix  $\mathbf{K}$  in (2) is built by considering a Mercer kernel  $K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  such that  $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ . An example kernel function that we considered throughout our work is the Gaussian kernel that for scale parameter  $\lambda$  is defined as  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\lambda^2}\right)$ .

Considering (2), we note that finding Kernel Minimum Enclosing balls boils down to finding optimal  $\boldsymbol{\mu}$ , which resides in the standard simplex  $\Delta^{n-1}$  and minimizes a convex function  $\mathcal{L}(\boldsymbol{\mu}) = \boldsymbol{\mu}^\top \mathbf{k} - \boldsymbol{\mu}^\top \mathbf{K} \boldsymbol{\mu}$ . Optimization settings of this kind can be easily solved iteratively using the Frank-Wolfe algorithm [3], which itself can be implemented as a recurrent neural network (see examples from [1, 2, 8, 10]). To this end, at each iteration  $t$ , the Frank-Wolfe algorithm evaluates the gradient of the negated dual Lagrangian  $\mathcal{L}(\boldsymbol{\mu})$  from (2), which amounts to  $\nabla \mathcal{L}(\boldsymbol{\mu}) = 2 \mathbf{K} \boldsymbol{\mu} - \mathbf{k}$ , and finds the vertex of  $\Delta^{n-1}$  for the update, that minimizes

$$\boldsymbol{\nu}_t = \underset{\mathbf{v}_j \in \mathbb{R}^n}{\operatorname{argmin}} \quad \mathbf{v}_j^\top [2 \mathbf{K} \boldsymbol{\mu}_t - \mathbf{k}] \approx \mathbf{g}_\beta (2 \mathbf{K} \boldsymbol{\mu}_t - \mathbf{k}), \tag{3}$$

where  $\boldsymbol{\nu}_t \in \mathbb{R}^n$  represent the current solution at  $t$ ,  $\mathbf{v}_j$  is the  $j$ th standard vector  $\mathbf{v}_j = [\delta_{j1}, \delta_{j2}, \dots, \delta_{jp}]^T$  (here  $\delta_{ji}$  represents the Kronecker delta) and, finally,

$\mathbf{g}_\beta(\mathbf{x})$  represents the soft-min operator. This operator is the smooth approximation of  $\text{argmin}_i$ , whose the  $i$ th entry defined as  $(\mathbf{g}_\beta(\mathbf{x}))_i = \frac{e^{-\beta x_i}}{\sum_j e^{-\beta x_j}}$  and has the limit

$$\lim_{\beta \rightarrow \infty} \mathbf{g}_\beta(\mathbf{x}) = \underset{\mathbf{v}_j \in \mathbb{R}^n}{\text{argmin}} \mathbf{v}_j^\top \mathbf{x} = \mathbf{v}_i. \quad (4)$$

Given that we can define the *convergent* iterative Frank-Wolfe updates [1] as

$$\boldsymbol{\mu}_{t+1} \leftarrow (1 - \eta_t) \boldsymbol{\mu}_t + \eta_t \mathbf{g}_\beta(2\mathbf{K}\boldsymbol{\mu}_t - \mathbf{k}), \quad (5)$$

where  $\eta_t \in [0, 1]$  is a monotonically decreasing step size. Rearranging the right-most expression in (5) as  $\mathbf{g}_\beta(2\mathbf{K}\boldsymbol{\mu}_t - \mathbf{k}) = \mathbf{g}_\beta(2\mathbf{K}\boldsymbol{\mu}_t + \bar{\mathbf{K}}\bar{\mathbf{1}})$ , where  $\bar{\mathbf{K}} = \text{diag}(\mathbf{k})$  and  $\bar{\mathbf{1}}$  is the vector of  $-1$ s defined as  $\bar{\mathbf{1}} = [-1, \dots, -1]^\top$ , allows us to interpret and implement these updates in terms of echo state networks [4]. That is, we can describe this machinery as a structurally constrained echo state network, in which we have the fixed input vector  $\bar{\mathbf{1}}$  containing  $-1$ s, the input weight matrix  $\bar{\mathbf{K}}$ ,  $n$  reservoir neurons with  $\mathbf{g}_\beta(\cdot)$  and  $2\mathbf{K}$  respectively being the nonlinear activation function and the reservoir weight matrices and  $\eta_t$  acting as a leaking rate for updating the Lagrange multipliers. Once optimal Lagrange multipliers have been found using the updates from (5), we can determine the kernelized radius and the squared magnitude of the center of the fitted ball  $\mathcal{B}$  respectively as  $r_* = \sqrt{\boldsymbol{\mu}_*^\top \mathbf{k} - \boldsymbol{\mu}_*^\top \mathbf{K} \boldsymbol{\mu}_*}$  and  $\mathbf{c}_*^\top \mathbf{c}_* = \boldsymbol{\mu}_*^\top \mathbf{K} \boldsymbol{\mu}_*$ , which will allow us to define a characteristic function defining the interior of  $\mathcal{B}$  [1]. Namely, using these equalities we can represent the inequality  $\|\mathbf{x} - \mathbf{c}_*\|^2 \leq r_*^2$  to check whether an arbitrary point  $\mathbf{x} \in \mathbb{R}^m$  within the ball  $\mathcal{B}$  by considering

$$f(\mathbf{x}) = \sqrt{K(\mathbf{x}, \mathbf{x}) - 2\bar{\mathbf{k}}^\top \boldsymbol{\mu}_* + \boldsymbol{\mu}_*^\top \mathbf{K} \boldsymbol{\mu}_*} - \sqrt{\boldsymbol{\mu}_*^\top \mathbf{k} - \boldsymbol{\mu}_*^\top \mathbf{K} \boldsymbol{\mu}_*}, \quad (6)$$

where  $\bar{\mathbf{k}} \in \mathbb{R}^n$  is defined as  $\bar{k}_i = K(\mathbf{x}, \mathbf{x}_i)$  [1]. That is,  $f(\mathbf{x}) > 0$  holds if  $\mathbf{x}$  is outside of the ball  $\mathcal{B}$ , whereas,  $f(\mathbf{x}) \leq 0$  is the case when  $\mathbf{x}$  is inside the ball  $\mathcal{B}$ . Though,  $f(\mathbf{x}) = 0$  only holds for the points with nonzero Lagrange multipliers that are the support vectors of  $\mathcal{B}$  and can be defined as  $\mathcal{S} = \{\mathbf{x}_i \mid \forall i \in [1, \dots, n] \wedge \mu_{i*} > 0\}$ . It is worth noting that, we can simplify (6) by grouping the  $l \leq n$  points in  $\mathcal{S}$  into a column data matrix  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_l] \in \mathbb{R}^{m \times l}$ , putting their corresponding multipliers in  $\boldsymbol{\sigma} \in \mathbb{R}^l$ , letting  $\mathbf{Q} \in \mathbb{R}^{l \times l}$  be the kernel matrix for the support vectors (i.e.  $Q_{ij} = K(\mathbf{s}_i, \mathbf{s}_j)$ ) and  $\mathbf{q} \in \mathbb{R}^l$  to contain its diagonal (i.e.  $\mathbf{q} = \text{diag}[\mathbf{Q}]$ ), which yields a simpler characteristic function

$$f(\mathbf{x}) = \sqrt{K(\mathbf{x}, \mathbf{x}) - 2\bar{\mathbf{k}}^\top \boldsymbol{\sigma} + \boldsymbol{\sigma}^\top \mathbf{Q} \boldsymbol{\sigma}} - \sqrt{\boldsymbol{\sigma}^\top \mathbf{q} - \boldsymbol{\sigma}^\top \mathbf{Q} \boldsymbol{\sigma}} \quad (7)$$

where as in (6),  $\bar{\mathbf{k}} \in \mathbb{R}^l$  is evaluated as  $\bar{k}_j = K(\mathbf{x}, \mathbf{s}_j)$  and we note that the term  $\sqrt{\boldsymbol{\sigma}^\top \mathbf{q} - \boldsymbol{\sigma}^\top \mathbf{Q} \boldsymbol{\sigma}}$  (which indeed amounts to  $r_*$ ) is does not depend on  $\mathbf{x}$ .

### 3 An Ensemble Approach for Novelty Discovery

Having explained how KMEBs are defined and can be computed so that we can determine their interior, we will now turn our attention to novelty discovery

by combining the characteristic functions of a set of balls. We note that, the characteristic function from (7) for a given ball  $\mathcal{B}$  can be used to label the points outside of the ball to be the novel points. In this case a query point  $\mathbf{x}$  is considered novel if  $f(\mathbf{x}) > 0$  and not novel for  $f(\mathbf{x}) \leq 0$ . Although this approach can capture novel points it might result in very restrictive or too general decision boundaries that respectively might result in detecting every query point to be novel or not novel (see Fig. 1d for the latter case). Both problems, however, can be avoided if we generalize this approach by combining the decisions of multiple balls. One approach for such a combination can be based on uniform voting [7]. That is, given a set of  $u$  KMEBs  $\mathcal{P} = \{\mathcal{B}_1, \dots, \mathcal{B}_u\}$ , that are trained considering a different setting, and  $f_i(\cdot)$  and  $\llbracket \cdot \rrbracket$  respectively indicating the characteristic function from (7) for ball  $\mathcal{B}_i$  and the Iverson bracket, we can assign the *novelty score* of a query point  $\mathbf{x}$  by evaluating  $z(\mathbf{x}) = \sum_{i=1}^u \llbracket f_i(\mathbf{x}) > 0 \rrbracket$  and, for instance, label  $\mathbf{x}$  to be novel if  $z(\mathbf{x}) \geq \lceil \frac{u}{2} \rceil$  (i.e.  $\mathbf{x}$  is outside of the majority of the balls in  $\mathcal{P}$  for an odd  $u$ ) and not novel if  $z(\mathbf{x}) < \lceil \frac{u}{2} \rceil$ . In the next section, we will empirically evaluate this methodology to detect novelty by showing two conceptual examples on benchmarking datasets.

**Table 1.** Novelty prediction results in terms of recall (**RC**), precision (**PR**), as well as the harmonic mean and geometric mean of both (respectively referred as **F1** and **GM**) for (a) the CBLC Face and (b) the MNIST datasets to respectively detect non-face images from face images and the images of digit 0 from the ones of 1. We benchmarked methods to detect novelty that consider the deviation from the sample mean (**MDEV**), matrix factorization (**MF**), euclidean MEBs (**EMEB**) and the ensemble of kernelized MEBs (**EKMEB**). The superior prediction results indicate that EKMEB can indeed be used for novelty discovery.

Method	RC	PR	F1	GM
<b>MDEV</b>	0.686	1.000	0.813	0.828
<b>MF</b>	0.711	0.999	0.831	0.843
<b>EMEB</b>	0.790	0.999	0.882	0.889
<b>EKMEB</b>	0.974	0.998	0.986	0.986

(a) MNIST dataset

Method	RC	PR	F1	GM
<b>MDEV</b>	0.043	1.000	0.082	0.207
<b>MF</b>	0.215	0.998	0.354	0.463
<b>EMEB</b>	0.095	1.000	0.173	0.308
<b>EKMEB</b>	1.000	0.949	0.974	0.974

(b) CBLC Faces

## 4 Empirical Results

We evaluated our method on the MNIST [5] and CBCL-face ([bit.ly/2Kw0VV6](http://bit.ly/2Kw0VV6)) datasets. For the former we trained models on the digit 1 aiming to obtain the 0s, whereas for the latter we leaned balls on faces to detect non-face novel images. So as to evaluate the precision of the detections, we divided the training data into 90/10 splits and the latter split is combined with the novel points, which resulted in training/evaluation datasets of cardinality values 6067/6598 and 2186/4791 for respectively the MNIST and CBCL-face datasets. We note

that for both examples, we constructed ensembles of KMEBs (i.e. distinct  $\mathcal{P}$  sets) with the Gaussian Kernel, whose scale values, in our case, were evenly spaced over specified intervals (as in Fig. 1) by considering  $u = 5$  KMEBs with  $\lambda$  ranging in  $[40, 60]$ . We also normalized the datasets to have zero mean and unit variance and always considered  $\beta = \infty$  for the softmin function (see (4)).

In Table 1, we compare our method against thresholding the tested points considering the maximum deviation from the sample mean vector [6], euclidean MEBs [1] (where we consider points outside of the ball as novel) and matrix factorization (MF) [8] based reconstruction to validate the use of kernel methods. For the first method, we label points in the test set as novel if the euclidean distance is larger than the furthest point to the sample mean. For the last method, we factorize the matrix with the number of latent factors  $k = 50$  using the alternating least squares method [9] and learn a threshold value based on the *worst* reconstruction error ( $l_2$ -norm). Unseen points with reconstruction error exceeding this threshold are considered novel. Table 1a and Table 1b respectively depict the prediction results for the MNIST and CBCL datasets, where we observe the superiority of ensemble KMEBs to detect novel datapoints.

## 5 Conclusion and Future Work

In this work, we introduced the idea of using ensemble of KEMBs for novelty discovery. We showed how we can construct ensembles of KEMBs and introduced a voting-based approach to detect novel data points. Our empirical evaluation yielded superior results over the use of mean deviation, euclidean MEBs and matrix factorization approaches. Our future work involves studying different ball selection as well as novelty determination strategies and extending the scope of the applications. Another line of future work is related to physical implementation of our methodology and in resource-constrained devices for applications in industrial domains such as for predictive maintenance.

## References

1. Bauckhage, C., Sifa, R., Dong, T.: Prototypes within minimum enclosing balls. In: Tetko, I.V., Kůrková, V., Karpov, P., Theis, F. (eds.) ICANN 2019. LNCS, vol. 11731, pp. 365–376. Springer, Cham (2019). [https://doi.org/10.1007/978-3-030-30493-5\\_36](https://doi.org/10.1007/978-3-030-30493-5_36)
2. Bauckhage, C.: A neural network implementation of Frank-Wolfe optimization. In: Lintas, A., Rovetta, S., Verschure, P.F.M.J., Villa, A.E.P. (eds.) ICANN 2017. LNCS, vol. 10613, pp. 219–226. Springer, Cham (2017). [https://doi.org/10.1007/978-3-319-68600-4\\_26](https://doi.org/10.1007/978-3-319-68600-4_26)
3. Frank, M., Wolfe, P.: An algorithm for quadratic programming. *Naval Res.* **3**, 95–110 (1956)
4. Jäger, H., Haas, H.: Harnessing nonlinearity: predicting chaotic systems and saving energy in wireless communication. *Science* **304**(5667), 78–80 (2004)
5. LeCun, Y., Boottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proc. IEEE* **86**(11), 2278–2324 (1998)

6. Pimentel, M.A., Clifton, D.A., Clifton, L., Tarassenko, L.: A review of novelty detection. *Sig. Process.* **99**, 215–249 (2014)
7. Rokach, L.: Ensemble methods for classifiers. In: Maimon, O., Rokach, L. (eds.) *Data mining and Knowledge Discovery Handbook*, pp. 957–980. Springer, Boston (2005). [https://doi.org/10.1007/0-387-25465-X\\_45](https://doi.org/10.1007/0-387-25465-X_45)
8. Sifa, R.: An overview of Frank-Wolfe optimization for stochasticity constrained interpretable matrix and tensor factorization. In: Kůrková, V., Manolopoulos, Y., Hammer, B., Iliadis, L., Maglogiannis, I. (eds.) *ICANN 2018. LNCS*, vol. 11140, pp. 369–379. Springer, Cham (2018). [https://doi.org/10.1007/978-3-030-01421-6\\_36](https://doi.org/10.1007/978-3-030-01421-6_36)
9. Sifa, R.: *Matrix and Tensor Factorization for Profiling Player Behavior*. LeanPub, British Columbia (2019)
10. Sifa, R., Paurat, D., Trabold, D., Bauckhage, C.: Simple recurrent neural networks for support vector machine training. In: Kůrková, V., Manolopoulos, Y., Hammer, B., Iliadis, L., Maglogiannis, I. (eds.) *ICANN 2018. LNCS*, vol. 11141, pp. 13–22. Springer, Cham (2018). [https://doi.org/10.1007/978-3-030-01424-7\\_2](https://doi.org/10.1007/978-3-030-01424-7_2)
11. Tax, D.M., Duin, R.P.: Data domain description using support vectors. In: *Proceedings of ESANN* (1999)