

# Semi-Supervised L2KC (S-L2KC) Classifier

Yue Lv

LY1251961956@163.COM

Jiangnan University, Wuxi, Jiangsu, China

**Editors:** Nianyin Zeng, Ram Bilas Pachori and Dongshu Wang

## Abstract

Building upon the density difference paradigm, a novel kernel classifier distinct from Support Vector Machines (SVM) — the L2-norm Kernel Classifier (L2KC) — has been developed. This methodology establishes an integrated squared error (ISE) criterion to estimate the true  $d_\gamma(x)$  through minimizing the L2-distance between  $d_\gamma(x)$  and  $\hat{d}_\gamma(x)$ , thereby achieving classification via explicit density difference representation. While L2KC demonstrates comparable accuracy to SVM with enhanced decision efficiency, its performance on real-world semi-supervised datasets requires improvement. To address this limitation, we propose the Semi-supervised L2KC (S-L2KC) by incorporating a locality-preserving projection (LPP) based manifold regularization term into the L2KC objective function. This integration effectively enforces the manifold assumption. Experimental results on benchmark datasets from the UCI and LIBSVM demonstrate that compared to L2KC, the proposed S-L2KC exhibits superior generalization capability, characterized by higher mean test accuracy with comparable or even smaller variance.

**Keywords:** L2KC; Semi-supervised Learning; Manifold Learning; Local Preserving Projections

## 1. Introduction

The L2-norm kernel classifier (L2KC) proposed by [Kim and Scott \(2010\)](#) is a classification algorithm analogous to Support Vector Machines (SVM) ([Cortes and Vapnik, 1995](#)). Consider a binary classification problem with samples and their corresponding class labels denoted as  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ , where  $\mathbf{x}_i \in R^d$  represents a d-dimensional sample and  $y_i \in \{1, -1\}$  its class label. Conventional kernel classifiers such as SVM employ a decision function (without offset term):

$$g(\mathbf{x}) = \text{sign} \left\{ \sum_{i=1}^n \alpha_i y_i k_\sigma(\mathbf{x}, \mathbf{x}_i) \right\} \quad (1)$$

In contrast, the L2KC proposed by [Kim and Scott \(2010\)](#) adopts a density difference paradigm. Let  $f_+(\mathbf{x})$  and  $f_-(\mathbf{x})$  denote the conditional probability density functions for the positive and negative class sample sets, respectively. The density difference is formulated as  $\hat{d}_\gamma(\mathbf{x}; \alpha) = f_+(\mathbf{x}) - \gamma f_-(\mathbf{x})$ . According to statistical decision theory, the optimal classifier takes the form:

$$g^*(\mathbf{x}) = \text{sign} \{ f_+(\mathbf{x}) - \gamma f_-(\mathbf{x}) \} \quad (2)$$

where  $\gamma$  denotes a user-defined fixed parameter that accounts for the class prior probabilities.

[Kim and Scott \(2010\)](#) redefine the class labels  $y_i \in \{-\gamma, 1\}$  and, based on the kernel density estimation (KDE) method, formulate the density functions for positive and negative class samples with weighting factors  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$  as:

$$\begin{aligned} \hat{f}_+(\mathbf{x}; \alpha) &= \sum_{i \in I_+} \alpha_i k_\sigma(\mathbf{x}, \mathbf{x}_i) \\ \hat{f}_-(\mathbf{x}; \alpha) &= \sum_{i \in I_-} \alpha_i k_\sigma(\mathbf{x}, \mathbf{x}_i) \end{aligned} \quad (3)$$

Then, an ISE function was formulated to estimate  $d_\gamma(\mathbf{x})$  via  $\hat{d}_\gamma(\mathbf{x}; \alpha) = \hat{f}_+(\mathbf{x}) - \gamma \hat{f}_-(\mathbf{x})$ , with parameter  $\alpha$  estimated by minimizing the L2 distance between  $\hat{d}_\gamma(\mathbf{x})$  and  $d_\gamma(\mathbf{x})$ . The mathematical expression is given by:

$$\begin{aligned} ISE(\alpha) &= \|\hat{d}_\gamma(\mathbf{x}; \alpha) - d_\gamma(\mathbf{x})\|_{L_2}^2 \\ &= \int \left( \hat{d}_\gamma(\mathbf{x}; \alpha) - d_\gamma(\mathbf{x}) \right)^2 d\mathbf{x} \\ &= \int \hat{d}_\gamma^2(\mathbf{x}; \alpha) d\mathbf{x} - 2 \int \hat{d}_\gamma(\mathbf{x}; \alpha) d_\gamma(\mathbf{x}) d\mathbf{x} + \int d_\gamma^2(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (4)$$

The optimization problem in Eq. (4) can be reformulated as a quadratic programming (QP) problem:

$$\hat{\alpha} = \arg \min_{\alpha \in A} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k_{\sqrt{2}\sigma}(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n c_i \alpha_i \quad (5)$$

The corresponding dual formulation is derived as:

$$\begin{aligned} \min_{w, \zeta_+, \zeta_-} J(w) &= \frac{1}{2} \mathbf{w}^2 + \zeta_+ + \zeta_- \\ s.t. \quad y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) &\geq c_i - \zeta_+, i \in I_+ \\ s.t. \quad y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) &\geq c_i - \zeta_-, i \in I_- \end{aligned} \quad (6)$$

where  $\phi_\sigma(\mathbf{x})$  denotes the kernel-induced mapping in the Gaussian reproducing kernel Hilbert space (RKHS). Differing from SVM, the L2-kernel classifier replaces the unit margin constraint with a flexible parameter  $c_i$ , implying that sufficiently large  $c_i$  values prioritize classification correctness of training samples  $\mathbf{x}_i$ . Furthermore, it employs a single smoothing factor  $\zeta_+$ ,  $\zeta_-$  per class without requiring sign-specific constraints.

## 2. S-L2KC

Currently, there are many semi-supervised learning algorithms, such as:

- Semi-Supervised Support Vector Machines (S3VMs) (Bennett and Demiriz, 1998).
- MixMatch, which combines consistency regularization and entropy minimization to improve pseudo-labeling (Berthelot et al., 2019).
- SimPLE, which optimizes classification networks using three training objectives and introduces pairwise loss for better feature learning (Hu and et al., 2021).
- UPS, which incorporates negative samples to filter out uncertain labels during training (Rizve et al., 2021).
- FreeMatch, which dynamically adjusts confidence thresholds based on the model's learning state (Wang and et al., 2022).

- PRCL, which enhances the Robustness of Contrastive Learning through Probabilistic Representation (Xie and et al., 2024).
- FixMatch, which integrates Consistency Regularization and Pseudo-Labeling Techniques (Sohn and et al., 2020).
- UniMatch, which extends FixMatch by strengthening the perturbation approach through unified image-level and feature-level perturbations (Yang et al., 2023).
- FlexMatch, which proposes a Curriculum Learning Approach – Curriculum Pseudo-Labeling (Zhang and et al., 2021).
- Consistent-Teacher: Triple Solutions for Stable Pseudo-labels (ASA + FAM-3D + GMM-Thresh) (Wang and et al., 2023).

To enhance the semi-supervised classification capability of L2KC, we incorporate a manifold regularization term based on locality-preserving projection (LPP) (He and Niyogi, 2005). The standard LPP objective function is formulated as:

$$\mathbf{w}^T \sum_{i=1}^n \sum_{j=1}^n (\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))(\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))^T \mathbf{w} \quad (7)$$

The Gram matrix  $(\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))(\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))^T$ , constructed from unlabeled instances, is incorporated as a manifold regularization term in the objective function to preserve the intrinsic geometric structure embedded in the data distribution.

By integrating the aforementioned manifold hypothesis into the optimization objective of S-L2KC through manifold regularization (Wang and et al., 2022), the enhanced objective function can be formulated as:

$$\begin{aligned} \min_{w, \zeta_+, \zeta_-} J(w) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{\mu}{2} \mathbf{w}^T A \mathbf{w} + \zeta_+ + \zeta_- \\ s.t. \quad y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) &\geq c_i - \zeta_+, i \in I_+ \\ s.t. \quad y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) &\geq c_i - \zeta_-, i \in I_- \end{aligned} \quad (8)$$

Formulating the Lagrangian function for the constrained optimization problem stated above leads to the following primal-dual relationship:

$$\begin{aligned} L(\mathbf{w}, \zeta_+, \zeta_-, \alpha, \beta) &= \hat{J}(w) \\ &- \sum_{i \in M, i \in I_+} \alpha_i (y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) - c_i + \zeta_+) \\ &- \sum_{i \in M, i \in I_-} \beta_i (y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) - c_i + \zeta_-) \end{aligned} \quad (9)$$

Through differentiation of the Lagrangian with respect to variables, the KKT optimality conditions enforce the following system of stationary equations:

$$\begin{aligned}
\frac{\partial L}{\partial w} = 0 &\rightarrow \mathbf{w} = \sum_{i \in M, i \in I_+} \alpha_i y_i (I + \mu A)^{-1} \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) + \sum_{i \in M, i \in I_-} \beta_i y_i (I + \mu A)^{-1} \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) \\
\frac{\partial L}{\partial \zeta_+} = 0 &\rightarrow \sum_{i \in M, i \in I_+} \alpha_i = 1 \\
\frac{\partial L}{\partial \zeta_-} = 0 &\rightarrow \sum_{i \in M, i \in I_-} \beta_i = 1 \\
\frac{\partial L}{\partial \alpha} = 0 &\rightarrow c_i - \zeta_+ = y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i) \\
\frac{\partial L}{\partial \beta} = 0 &\rightarrow c_i - \zeta_- = y_i \mathbf{w}^T \phi_{\sqrt{2}\sigma}(\mathbf{x}_i)
\end{aligned} \tag{10}$$

Through algebraic manipulation and substitution of the stationary conditions, the optimization problem reduces to the following canonical form:

$$\begin{aligned}
B = \{ \alpha \mid \sum_{i \in I_+} \alpha_i = \sum_{i \in I_-} \alpha_i = 1, i \in N \} \\
\frac{1}{2} \sum_{i \in N} \sum_{j \in M} \alpha_i \alpha_j y_i y_j \phi_{\sqrt{2}\sigma}(\mathbf{x}_i)^T (I + \mu A)^{-1} \phi_{\sqrt{2}\sigma}(\mathbf{x}_j) = \sum_{i \in N} \alpha_i c_i - \zeta
\end{aligned} \tag{11}$$

The primal optimization problem is thus reduced to a convex quadratic programming problem, formally expressed as:

$$\hat{\alpha} = \arg \min_{\alpha \in B} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^M \alpha_i \alpha_j y_i y_j \phi_{\sqrt{2}\sigma}(\mathbf{x}_i)^T (I + \mu A)^{-1} \phi_{\sqrt{2}\sigma}(\mathbf{x}_j) - \sum_{i=1}^n c_i \alpha_i \tag{12}$$

Solving the quadratic programming problem yields the dual coefficients, which can be substituted into the decision function to obtain the final S-L2KC classifier formulation:

$$g^*(\mathbf{x}) = \text{sign} \left\{ \sum_{i=1}^N \alpha_i y_i \phi_{\sqrt{2}\sigma}(\mathbf{x}_i)^T (I + \mu A)^{-1} \phi_{\sqrt{2}\sigma}(\mathbf{x}) \right\} \tag{13}$$

### 3. Experimental results and analysis

This paper utilizes several datasets from the UCI and LIBSVM databases, with the sample sizes and feature information of the relevant datasets detailed in Table 1. The study employs a 10-folds cross-validation approach, where the preprocessed dataset is divided into ten parts. Nine parts are selected as the training set, and the remaining one serves as the test set. To more accurately demonstrate the classification performance of the algorithm, we use the average accuracy and variance obtained from performing 50 iterations of ten-fold cross-validation on different datasets as the performance metrics.

Table 1: Dataset details

Dataset	Sample Size	Feature Count
<i>Diabetes</i>	768	8
<i>German</i>	1000	20
<i>Heart</i>	270	13
<i>Image</i>	2310	18
<i>Ionosphere</i>	351	34
<i>Breast-cancer</i>	277	9
<i>sonar</i>	208	60
<i>Waveform</i>	5000	21
<i>Splice</i>	3175	60
<i>Ringnorm</i>	7400	20

Based on real-world experience, the proportion of unlabeled samples in a dataset is generally high. However, considering extreme conditions is not conducive to the algorithm’s learning in semi-supervised scenarios. Therefore, in the experimental design, this paper sets the ratio ( $r$ ) of unlabeled samples to total samples as (0.6, 0.7, 0.8) to test the algorithm’s performance under reasonable proportions of unlabeled data and examines the algorithm’s robustness by varying this ratio. Table 2 presents the relevant parameter settings for this experiment. During the experiment, the hyperparameters of the L2KC and S-L2KC algorithms were determined through grid search. A Gaussian kernel function was employed, with the hyperparameter ( $\sigma$ ) ranging from  $10^2$  to  $10^{-2}$  to avoid kernel matrix sparsity caused by excessively small values, while allowing  $\sigma$  to be increased when unlabeled samples dominate for enhanced global stability.  $\mu$  controls the strength of the manifold regularization term, with a range from  $10^2$  to  $10^{-2}$ . When  $\mu$  is too small, it results in weak manifold constraints, making the model approach supervised learning. Conversely, when  $\mu$  is too large, it enforces overly rigid geometric structure preservation.

Table 2: Dataset details

Parameter	Parameter Value
$r$	0.6, 0.7, 0.8
$\sigma$	$10^{-2}$ to $10^2$
$\mu$	$10^{-2}$ to $10^2$

Tables 3, 4, and 5 present the training accuracy (expressed as mean  $\pm$  standard deviation) of the L2KC and S-L2KC algorithms on the aforementioned UCI and LIBSVM datasets. In these tables, values where the S-L2KC algorithm achieves an average accuracy improvement greater than 0.6% are bolded, and values where the variance reduction exceeds 0.2% are underlined.

The experimental results demonstrate that S-L2KC generally achieves higher test accuracy and lower test variance compared to L2KC, with the improvement becoming more pronounced as the proportion of unlabeled samples increases (from 60% to 70%). However, when the proportion of unlabeled samples further increases, the performance gain of S-L2KC diminishes slightly, though it still maintains an advantage over L2KC. Additionally, the CPU time of both algorithms remains very close (less than 1ms).

Table 3: Average test accuracies and standard deviations ( $r = 0.6$ )

Dataset	L2KC	S-L2KC
<i>Diabetes</i>	72.725 $\pm$ 2.374	<b>73.348<math>\pm</math>2.301</b>
<i>German</i>	70.104 $\pm$ 1.951	70.483 $\pm$ 1.733
<i>Heart</i>	82.152 $\pm$ 4.137	82.471 $\pm$ 4.109
<i>Image</i>	70.972 $\pm$ 6.433	71.318 $\pm$ 6.201
<i>Ionosphere</i>	70.011 $\pm$ 3.730	70.425 $\pm$ 3.703
<i>Breast-cancer</i>	73.384 $\pm$ 4.342	73.377 $\pm$ 4.293
<i>sonar</i>	81.949 $\pm$ 3.551	<b>83.072<math>\pm</math>3.216</b>
<i>Waveform</i>	84.573 $\pm$ 0.733	<b>85.541<math>\pm</math>0.701</b>
<i>Splice</i>	59.443 $\pm$ 4.307	<b>60.172<math>\pm</math>4.103</b>
<i>Ringnorm</i>	97.479 $\pm$ 0.239	97.470 $\pm$ 0.225

Table 4: Average test accuracies and standard deviations ( $r = 0.7$ )

Dataset	L2KC	S-L2KC
<i>Diabetes</i>	72.304 $\pm$ 2.401	<b>73.222<math>\pm</math>2.189</b>
<i>German</i>	69.696 $\pm$ 1.961	<b>70.384<math>\pm</math>1.863</b>
<i>Heart</i>	81.183 $\pm$ 4.132	81.506 $\pm$ 3.891
<i>Image</i>	70.505 $\pm$ 6.403	<b>71.248<math>\pm</math>6.199</b>
<i>Ionosphere</i>	69.991 $\pm$ 3.730	70.375 $\pm$ 3.582
<i>Breast-cancer</i>	73.408 $\pm$ 4.331	73.414 $\pm$ 4.248
<i>sonar</i>	82.079 $\pm$ 3.527	<b>82.971<math>\pm</math>3.321</b>
<i>Waveform</i>	84.907 $\pm$ 0.744	85.319 $\pm$ 0.729
<i>Splice</i>	59.731 $\pm$ 4.291	<b>61.231<math>\pm</math>4.031</b>
<i>Ringnorm</i>	97.239 $\pm$ 0.241	97.235 $\pm$ 0.243

Table 5: Average test accuracies and standard deviations ( $r = 0.8$ )

Dataset	L2KC	S-L2KC
<i>Diabetes</i>	72.337 $\pm$ 2.351	<b>72.941<math>\pm</math>2.348</b>
<i>German</i>	69.005 $\pm$ 1.922	69.283 $\pm$ 1.917
<i>Heart</i>	81.037 $\pm$ 4.107	<b>81.701<math>\pm</math>4.008</b>
<i>Image</i>	70.672 $\pm$ 6.389	<b>71.342<math>\pm</math>6.388</b>
<i>Ionosphere</i>	69.759 $\pm$ 3.691	70.048 $\pm$ 3.696
<i>Breast-cancer</i>	73.384 $\pm$ 4.342	73.847 $\pm$ 4.339
<i>sonar</i>	81.992 $\pm$ 3.545	<b>83.211<math>\pm</math>3.541</b>
<i>Waveform</i>	84.986 $\pm$ 0.740	84.887 $\pm$ 0.741
<i>Splice</i>	62.358 $\pm$ 4.288	<b>63.517<math>\pm</math>4.207</b>
<i>Ringnorm</i>	97.171 $\pm$ 0.244	97.159 $\pm$ 0.240

## 4. Conclusion

This paper presents an introduction to L2KC while examining its theoretical foundations and unresolved challenges in semi-supervised learning scenarios. To enhance L2KC’s capability for semi-

supervised learning, we propose a manifold regularization term based on Locality Preserving Projections (LPP) in Gram matrix form. By incorporating this regularization term into the original L2KC framework, we develop an improved algorithm termed S-L2KC, complete with theoretical analysis and predictive functions. A series of experiments are then conducted to evaluate its semi-supervised learning performance. The experimental results demonstrate that on selected UCI and LIBSVM datasets, S-L2KC achieves either superior test accuracy or comparable accuracy with reduced variance compared to L2KC, indicating enhanced generalization capability. Detailed analysis reveals that the performance improvement becomes more significant when the proportion of unlabeled samples increases from 60% to 70%. While further increasing the unlabeled sample ratio to 80% leads to slightly diminished effects, S-L2KC maintains consistent advantages in both accuracy and variance metrics. These findings confirm that S-L2KC effectively leverages information from unlabeled samples to improve classification performance and generalization ability in semi-supervised learning scenarios when compared to the original L2KC approach.

However, the performance of S-L2KC is highly dependent on the adaptability of the selected kernel function. If the kernel function fails to effectively capture the underlying structure of the data, the model’s performance may degrade significantly. Additionally, in scenarios involving high noise or sparse data, the manifold assumption may not hold.

## References

- Kristin P. Bennett and Ayhan Demiriz. Semi-supervised support vector machines. In *Advances in Neural Information Processing Systems*, pages 368–374, 1998. doi: 10.1007/978-0-387-77501-2\_7.
- David Berthelot, Nicholas Carlini, and et al. Mixmatch: A holistic approach to semi-supervised learning. *arXiv preprint arXiv:1905.02249*, 2019. doi: 10.48550/arXiv.1905.02249.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995. doi: 10.1023/A:1022627411411.
- Xiaofei He and Partha Niyogi. Locality preserving projections. *Neural Information Processing Systems*, pages 153–160, 2005. doi: 10.1007/978-3-319-73830-7\_13.
- Zijian Hu and et al. Simple: Similar pseudo label exploitation for semi-supervised classification. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10428–10437, 2021. doi: 10.48550/arXiv.2103.16725.
- Joo Seuk Kim and C. Scott.  $l_2$  kernel classification. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 32(10):1822–1831, 2010. doi: 10.1109/TPAMI.2009.188.
- Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. *arXiv preprint arXiv:2101.06329*, 2021. doi: 10.48550/arXiv.2101.06329.
- Kihyuk Sohn and et al. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv preprint arXiv:2001.07685*, 2020. doi: 10.48550/arXiv.2001.07685.

- X. Wang and et al. Consistent-teacher: Towards reducing inconsistent pseudo-targets in semi-supervised object detection. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3240–3249, 2023. doi: 10.48550/arXiv.2209.01589.
- Yidong Wang and et al. Freematch: Self-adaptive thresholding for semi-supervised learning. *arXiv preprint arXiv:2205.07246*, 2022. doi: 10.48550/arXiv.2205.07246.
- Haoyu Xie and et al. Prcl: Probabilistic representation contrastive learning for semi-supervised semantic segmentation. *International Journal of Computer Vision*, 132(10):2495–2514, 2024. doi: 10.1007/s11263-024-02016-8.
- L. Yang, L. Qi, L. Feng, and et al. Revisiting weak-to-strong consistency in semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7236–7246, 2023. doi: 10.48550/arXiv.2208.09910.
- Bowen Zhang and et al. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *arXiv preprint arXiv:2110.08263*, 2021. doi: 10.48550/arXiv.2110.08263.