# Universiteit Leiden
## The Netherlands

# Application of functional kernel hypothesis testing for channel selection in time series classification
Huang, Q.; Bäck, T.H.W.; Stein, N. van

# Application of Functional Kernel Hypothesis Testing for Channel Selection in Time Series Classification

Qi Huang
*LIACS, Leiden University*
Leiden, The Netherlands
0009-0007-4989-135X ⦿

Thomas Bäck
*LIACS, Leiden University*
Leiden, The Netherlands
0000-0001-6768-1478 ⦿

Niki van Stein
*LIACS, Leiden University*
Leiden, The Netherlands
0000-0002-0013-7969 ⦿

*Abstract*—**Multi-variate time series classification tasks are prevalent in various real-world engineering domains, including but not limited to activity recognition and anomaly detection. However, due to the abundance of sensors available, selecting the appropriate channels for successful classification can be a daunting task. In this study, we propose to use a two-sample hypothesis test, to determine the relevance of channels in time series classification tasks. We illustrate the industrial usecase that motivated this algorithm and validate our approach on open-source benchmarks. The proposed method has the potential to address the challenge of channel selection in multi-variate time series classification tasks and can significantly impact various real-world engineering applications.**

*Index Terms*—**Channel Selection, Feature Selection, Time Series Classification**

## I. INTRODUCTION

Let $X \sim \mathbb{P}_X$ be a stochastic process (channel), where $x_t = X(t)$ is its realization on a time point $t$ and we let $X_I = [x_1, x_2, \ldots, x_n]$ to represent the segment of $X$ observed over a time interval $I$ with $n$ measurement points ($t_i$). Additionally, we define time series $T$, a segment-collection of $X$, as $T = \{T_1, T_2, \ldots, T_k\}$, where $T_j = X_{I_j}$. Time series classification requires building a machine learning model to predict the class labels of unseen time series based on $k$ historical data-label observations $T$ and labels $L$. Assuming for each time interval $I$, we observe $M \geq 2$ stochastic processes simultaneously, i.e., $T_j = (X_{1,I_j}, \ldots, X_{M,I_j})$, where $X_{i,I_j}$ is the segment of channel $i$ on interval $I_j$. Then the problem becomes multivariate or multi-channel time series classification.

When dealing with classification tasks in complex industrial environments, one of the challenges is the large number of sensors involved. Each sensor generates a sequence of outputs (a channel). It is often unclear which sensors have a direct or indirect relation with the classification target, posing difficulties in selecting the relevant channels for classification. To provide a concrete example, we consider the industrial use-case of classifying the activity of the world's largest crane vessel, the Sleipnir, owned by the Heerema company. This vessel is equipped with over $m = 6000$ sensors, measuring various subsystems such as thrusters, ballast pumps, cranes, and engines. Depending on the activity to be classified,

many of these channels are potentially relevant, and selecting a feasible subset of channels is a challenging and time-consuming task. Hence, this research aims to handle the following scenario: **Given a TSC task on time series** $T$ **with** $M$ **channels (sensors), find the most** $P \leq M$ **relevant channels for predicting the segment-wise labels** $L$. For a robust TSC classifier $f$, a train set $D_{tr}$, a test set $D_{te}$, and a performance metric $P(D_{te})$ where a larger value means better performance. Herewith, we define a channel $X_{i,*}$ is *truly relevant* if it satisfies the following criteria:

- Fit $f$ exclusively on $X_{i,train}, L$ and evaluate performance on the test data $D_{i,test}$. Then $P(D_{te}) - P(D_{i,te}) \leq \epsilon$.
- Remove $X_{i,*}$ from $D_{tr}$ and $D_{te}$, and let the new train and test sets be $D_{tr}^i$ and $D_{te}^i$. Train $f$ on $D_{tr}^i$ and obtain its performance $P(D_{te}^i)$, then $P(D_{te}) - P(D_{te}^i) \geq \delta$,

where both $\epsilon$ and $\delta \geq 0$ are scenario-dependent thresholds.

It is possible to determine the relevance by exhaustively doing leave-one-channel-out runs. However, this violates the principle of channel selection, namely, saving computational costs. Thus, we propose a methodology that can quantify and estimate this relevance based on a functional kernel-based two-sample test [1] to determine the relevance of channels for a given TSC task. Due to the confidentiality of our industrial data, we verify our proposed method on several open-source benchmarks and demonstrate the effectiveness and limitations of the proposed approach.

## II. RELATED WORK

Although various feature selection methods exist for classical tabular machine learning, there is a notable research gap in terms of selecting channels for multi-channel time series tasks. Some existing methods employ static feature extraction to determine the relevant channels and compute inner-correlations, such as entropy, to rank the performance of different learning algorithms on different feature subsets or employ a hybrid of these two methods [2], [3]. Another approach is to study the correlations between channels, including joint correlations, using Pearson's correlation [4], [5]. However, we argue that Pearson's correlation does not account for the time-dependency of channels, which makes it imperfect for time-series analysis.

## III. METHODOLOGY

*Time space to function space:* It is assumed that each univariate time series is a function supported over a vector space, and can be represented by a finite vector [6].

*Functional kernel hypothesis test:* A functional kernel two-sample test aims to test the distributional equivalence between two functions supported by the same vector space. The core idea is to use a reproducing kernel defined on functional space to construct and test the metric discrepancies between the mean embedding of two functions [1].

Based on these preliminaries, for given $k$ time intervals $I = \{I_1, \ldots, I_k\}$, we observe a univariate time series channel $X = \{X_{I_1}, \ldots, X_{I_k}\}$ and its element-wise labels $L = \{L_1, \ldots, L_k\}$, i.e., pairs of $(X_{I_j}, L_j)$. Additionally, we assume the time difference (start to end) for all time intervals is the same but the sampling strategy for measurement points in each $I_j$ can be different. Consequently, it is feasible for us to consider that all $X_{I_j}$ are defined in the same space but with different observation grids. Moreover, let the functional kernel two-sample test algorithm be $g(F_p, F_q; \kappa)$, where $F_p \sim \mathbb{P}$ and $F_q \sim \mathbb{Q}$ are functional distributions and $\kappa$ is the functional reproducing kernel. Then the $g$-based relevance of $X$ w.r.t $L$ is determined as follows:

1) Suppose there are $m$ unique labels in $L$, namely, $A = \{A_1, \ldots, A_m\}$ and note that $X \to A$ is surjective-only. For each unique $A_i \in A$, we find the two subsets $X_i = \{X_{I_j} \mid L_j = A_i, \ \forall j \in [1, k]\}$ and $X_i^* = \{X \setminus X_i\}$.

2) For each pair of $(X_i, X_i^*)$, we determine their two-sample-test-based test power, namely, $t_i = g(X_i, X_i^*; \kappa)$. It is obvious that here we consider $X_i$ and $X_i^*$ to be two different functions supported by $I$.

3) The relevance score of $X$ w.r.t $L$ bounded by $\kappa$, namely, $R(X \mid L, \kappa)$ can then be written as:

$$R(X \mid L, \kappa) = \sum_{i=1}^{m} \frac{t_i}{m}$$

4) Suppose we have $C$ valid functional reproducing kernels, i.e., $\Omega = \{\kappa_1, \ldots, \kappa_C\}$. We do procedures 1 to 3 for $C$ times, each time using a unique $\kappa_j \in \Omega$. Then the final channel relevance of $X$ w.r.t $L$ is defined as:

$$R(X \mid L) = \max_{j=1}^{C} R(X \mid L, \kappa_j)$$

It is beneficial to use multiple kernels since each kernel is limited by its form, and can only capture certain types of hypothesis classes of the functions[1].

## IV. EXPERIMENTS

*Setup:* Recall the initial goal of finding the top $P \leq M$ relevant channels for TSC task. For efficiency, we generalize *relevance* to multi-channels by considering or removing a batch of channels that fall into a certain range of relevance. To verify the soundness of our proposed channel relevance,

---

[1]However, considering two channels both with high relevancy but are based on fundamentally different kernels, then multiple kernels can be misleading since a classifier may fail to process two intrinsically different function classes.

---

we perform benchmarking on the well-established UCR time series suites [7]. The experiments are organized as follows:

1) Given a $m$-channels multi-variate TSC dataset ($T = \{X_{1,*}, \ldots, X_{m,*}\}$). For each channel $X_{i,*}$, we compute its channel relevance score $R(X_{i,*}|L)$ based on the aforementioned algorithm.

2) Grouping channels by their channel relevance, e.g., $S_{[a,b]} = \{X_{j,*} \mid R(X_{j,*}|L) \in [a,b], \ \forall j \in [1, \ldots, m]\}$.

3) Training and testing the performance of a classifier on each of the $S_{[a,b]}$.

For each of the used datasets, the classifier which yields the best results w.r.t the extensive review by Ruiz et al. [8] is used for benchmarking, where each classifier has been independently run 10 times on the dataset.

In addition to the benchmark settings, we outline the five kernels that are utilized in this experiment. Suppose we are distinguishing between two functional distributions $F_p \sim \mathbb{P}$ and $F_q \sim \mathbb{Q}$ supported over $\mathbb{R}^d$, the kernel can be written as:

$$\kappa(F_p, F_q) = e^{-\frac{\|T(F_p) - T(F_q)\|^2}{2\eta^2}},$$

where $\eta$ is intrinsically the bandwidth of the Gaussian kernel, $T$ is a Borel measurable, continuous, and injective transformation that is capable of mapping both functions into a real and separable Hilbert space. In this study, we consider the following five kernels from [1]:

- Standard: $T(a) = a$, the original Gaussian RBF kernel
- Cosine (Cos): point-wise cosine transformation
- Squaring feature expansion (Sqr): summing up two kernels. One with $T(a) = a$ and another one with $T(a) = a^2$
- FPCA: $T$ is functional PCA.
- Covariance (Cov): it is different from other kernels, the outcome is the inner product between $F_p$ and $F_q$.

*Results:* Four symbolic results obtained on four datasets from the UCR multivariate suite [9] are shown in Table I,II,III and IV.

TABLE I
RESULTS OBTAINED ON *RacketSports*

|  | Range of channel relevance | | |
|---|---|---|---|
|  | (0.4,0.6] | (0.6,0.8] | (0.8,1.0] |
| F1 | 0.80 | 0.82 | 0.83 |
|  | 0.85 | | 0.83 |
|  | 0.80 | **0.88** | |
|  | | 0.86 | |
| N | 2 | 2 | 2 |

TABLE II
RESULTS OBTAINED ON *HandMovementDirection*

|  | Range of channel relevance | | | |
|---|---|---|---|---|
|  | [0,0.1] | (0.1,0.2] | (0.2,0.3] | (0.3,0.4] |
| F1 | 0.39 | 0.45 | 0.46 | 0.48 |
|  | 0.43 | | **0.55** | |
|  | 0.47 | | | 0.48 |
|  | 0.39 | | 0.52 | |
|  | | 0.52 | | |
| N | 4 | 3 | 2 | 1 |

We first give an example of how to read the tables, "F1" means the (weighted) $f1$ scores (averaged over 10 independent runs) and $N$ means the number of channels that falls into the relevance range (denoted at the top of each column). The two results in the second row of Table II shows the $f1$ scores $0.43$ and $0.55$ obtained by the classifier on channels whose relevance are between $[0, 0.2]$ and $(0.2, 0.4]$, respectively. The darker the color of a cell is, the better the performance. The best results per data set are shown in **bold** font.

It can be seen from Table I, by leaving out the two channels with *lower relevance*, the classifier yields better performance than using all channels. Similar and more encouraging results can be seen in Table II, where only relying on the top 3 relevant channels outperforms using all 10 channels.

### TABLE III
#### RESULTS OBTAINED ON *Heartbeat*

| | Range of channel relevance | | | | | | |
|---|---|---|---|---|---|---|---|
| | [0,0.15] | (0.15,0.3] | [0.3,0.45] | [0.45,0.6] | (0.6,0.75] | (0.75,0.9] | (0.9,1.0] |
| | 0.70 | 0.71 | 0.70 | 0.73 | 0.76 | 0.69 | 0.76 |
| | 0.70 | | | | **0.77** | | |
| | | 0.71 | | | **0.77** | | |
| F1 | | 0.70 | | | | 0.76 | |
| | | | 0.72 | | | **0.77** | |
| | | | 0.74 | | | | 0.76 |
| | | | | 0.74 | | | 0.76 |
| | | | | **0.77** | | | |
| N | 1 | 2 | 2 | 5 | 10 | 10 | 31 |

### TABLE IV
#### RESULTS OBTAINED ON *NATOPS*

| | Range of channel relevance | | | | | | |
|---|---|---|---|---|---|---|---|
| | (0.3,0.4] | (0.4,0.5] | (0.5,0.6] | (0.6,0.7] | (0.7,0.8] | (0.8,0.9] | (0.9,1.0] |
| | 0.69 | 0.67 | 0.81 | 0.78 | 0.75 | 0.60 | 0.88 |
| | 0.69 | | | **0.91** | | | |
| | | 0.75 | | | **0.91** | | |
| F1 | | 0.82 | | | **0.91** | | |
| | | 0.87 | | | | 0.89 | |
| | | | 0.89 | | | | 0.87 |
| | | | 0.91 | | | | 0.88 |
| | | | | 0.91 | | | |
| N | 1 | 4 | 4 | 3 | 2 | 2 | 8 |

In Table III, by leaving out 10 channels the classifier is still performing as good as including all channels. One could even opt to leave out the 30 least informative channels and only lose $0.01$ on the $f1$ score. In principle, the values on the right side should be higher than the ones on the left side, as these include the most relevant channels, but this is not guaranteed since the proposed channel relevance is based on correlation instead of causality and the inter-channel information shall be taken into account. From Table IV one can observe that leaving out the $8$ most relevant channels leads to the same $f1$ score as leaving out the $9$ least relevant channels. Looking at the channel relevance scores of the NATOPS benchmark, we can observe that actually all channels show some relevance and therefore this does make sense.

Lastly, Table V shows the number of times the relevance of a channel is determined by a certain type of kernel. Looking at the **bold** numbers, it is not hard to see that the squaring function expansion is the most confident kernel on three out of four problems (Heartbeat (HB), RacketSports (RS), and NATOPS), which suggests the original signals could possess strong polynomial inner-correlations. On the contrary, it can be

assumed that there exist more sine-cosine correlations within the data of HandMovementDirection (HMD).

### TABLE V
#### NUMBER OF TIMES THAT A KERNEL OBTAINS THE HIGHEST RELEVANCE

| | Standard | Cos | Sqr | Cov | FPCA | Total |
|---|---|---|---|---|---|---|
| HMD | 1 | **7** | 0 | 0 | 2 | 10 |
| HB | 0 | 0 | **61** | 0 | 0 | 61 |
| RS | 0 | 0 | **4** | 2 | 0 | 6 |
| NATOPS | 1 | 1 | **19** | 3 | 0 | 24 |

## V. CONCLUSIONS AND OUTLOOK

A novel functional kernel-based two-sample test approach is proposed to determine the channel relevance in multivariate time series classification tasks. The effectiveness of the approach is validated on a wide set of benchmark TSC problems. We show that the two-sample test methodology, combining different static kernels, works very well in order to reduce the number of channels and even increase classification accuracy (f1 score) in some cases. Although the experimental results show that one kernel is promising most of the time, choosing the best static functional kernel as well as selecting hyperparameters can however still pose a challenge for real-world applications, and the inter-channel relevance is not taken into account using the proposed approach. For future research directions, a learned kernel could be interesting to solve the kernel selection problem.

## REFERENCES

[1] G. Wynne and A. B. Duncan, "A Kernel Two-Sample Test for Functional Data," *Journal of Machine Learning Research*, vol. 23, no. 73, pp. 1–51, 2022. [Online]. Available: http://jmlr.org/papers/v23/20-1180.html

[2] T. Alotaiby, F. E. A. El-Samie, S. A. Alshebeili, and I. Ahmad, "A review of channel selection algorithms for EEG signal processing," *EURASIP Journal on Advances in Signal Processing*, vol. 2015, no. 1, p. 66, Aug. 2015. [Online]. Available: https://doi.org/10.1186/s13634-015-0251-9

[3] M. Z. Baig, N. Aslam, and H. P. H. Shum, "Filtering techniques for channel selection in motor imagery EEG applications: a survey," *Artificial Intelligence Review*, vol. 53, no. 2, pp. 1207–1232, Feb. 2020. [Online]. Available: https://doi.org/10.1007/s10462-019-09694-8

[4] J. Jin, Y. Miao, I. Daly, C. Zuo, D. Hu, and A. Cichocki, "Correlation-based channel selection and regularized feature optimization for MI-based BCI," *Neural Networks*, vol. 118, pp. 262–270, Oct. 2019. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0893608019301960

[5] D. D. Chakladar and S. Chakraborty, "EEG based emotion classification using "Correlation Based Subset Selection"," *Biologically Inspired Cognitive Architectures*, vol. 24, pp. 98–106, Apr. 2018. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2212683X18300227

[6] L. Horváth and P. Kokoszka, "Inference for Functional Data with Applications," ser. Springer Series in Statistics, vol. 200. New York, NY: Springer New York, 2012. [Online]. Available: https://link.springer.com/10.1007/978-1-4614-3655-3

[7] H. A. Dau, A. Bagnall, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A. Ratanamahatana, and E. Keogh, "The UCR time series archive," *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 6, pp. 1293–1305, Nov. 2019.

[8] A. P. Ruiz, M. Flynn, J. Large, M. Middlehurst, and A. Bagnall, "The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances," *Data Mining and Knowledge Discovery*, vol. 35, no. 2, pp. 401–449, Mar. 2021. [Online]. Available: https://doi.org/10.1007/s10618-020-00727-3

[9] H. A. Dau, E. Keogh, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A. Ratanamahatana, Yanping, B. Hu, N. Begum, A. Bagnall, A. Mueen, G. Batista, and Hexagon-ML, "The ucr time series classification archive," October 2018, https://www.cs.ucr.edu/~eamonn/time_series_data_2018/.