# Kernelized Sorting

Authors

**Novi Quadrianto**[1] | **Le Song**[2] | **Alex J. Smola**[3]

## 1: RSISE, ANU & SML, NICTA | 2: SCS, CMU | 3: Yahoo! Research

## Abstract

Matching pairs of objects from different domains is a fundamental operation in data analysis. It typically requires the definition of a similarity measure between the classes of objects to be matched. For many cases, we may be able to design a cross domain similarity measure based on prior knowledge or to observe one based on the co-occurence of such objects. In some cases, however, such a measure may not exist or it may not be given to us beforehand.
We develop an approach which is able to perform matching by requiring a similarity measure only within each of the classes. This is achieved by maximizing the dependency between matched pairs of observations by means of the Hilbert-Schmidt Independence Criterion. This problem can be cast as one of maximizing a quadratic assignment problem with special structure and we present a simple algorithm for finding a locally optimal solution.

## Problem Statement

Assume we are given two collections of documents purportedly covering the same content, written in two different languages. Can we determine the correspondence between these two sets of documents without using a dictionary?
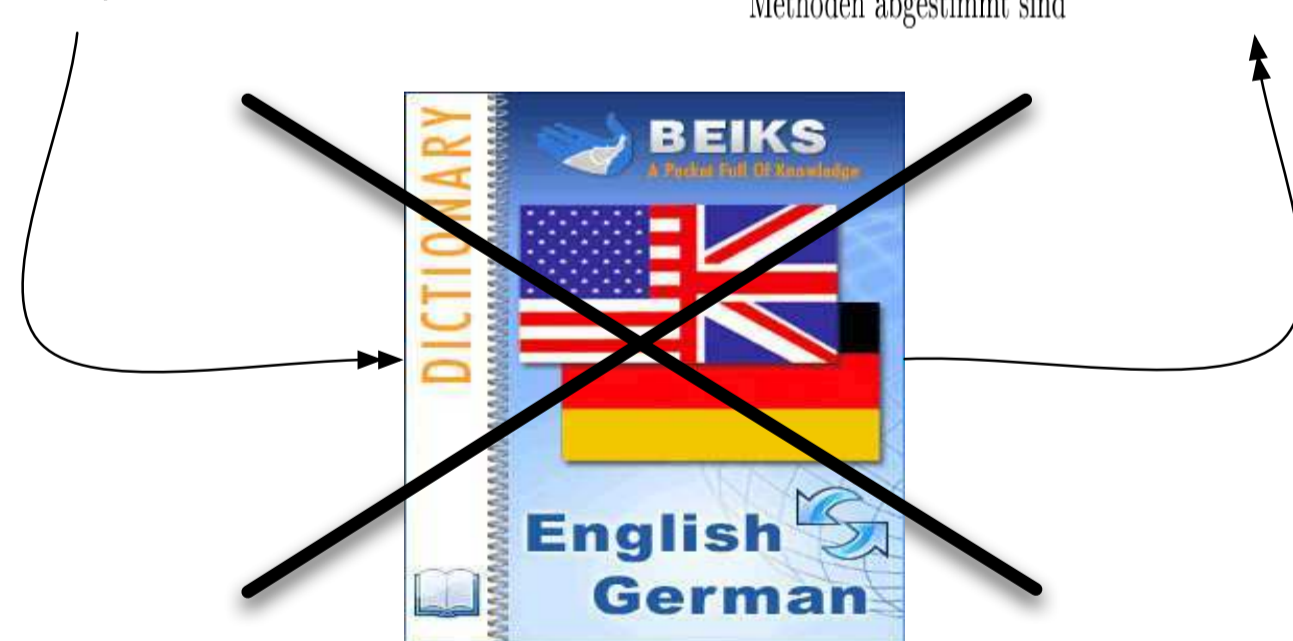
ENGLISH

Support Vector (SV) Machines combine several techniques from statistics, machine learning and neural networks. One of the most important ingredients are kernels, i.e. the concept of transforming linear algorithms into nonlinear ones via a map into feature spaces. The present work focuses on the following issues:

- Extensions of Support Vector Machines.
- Extensions of kernel methods to other algorithms such as unsupervised learning.
- Capacity bounds which are particularly well suited for kernel methods.

GERMAN

Support Vektor (SV) Maschinen verbinden verschiedene Techniken der Statistik, des maschinellen Lernens und Neuronaler Netze. Eine Schlüsselposition fällt den Kernen zu, d.h. dem Konzept, lineare Algorithmen durch eine Abbildung in Merkmalsräume nichtlinear zu machen. Die Dissertation behandelt folgende Aspekte:

- Erweiterungen des Support Vektor Algorithmus
- Erweiterungen und Anwendungen kernbasierter Methoden auf andere Algorithmen wie das unüberwachte Lernen
- Abschätzungen zur Generalisierungsfähigkeit, die besonders auf kernbasierte Methoden abgestimmt sind

(Formal) problem formulation:
Given

- two sets of observations $X = \{x_1, \ldots, x_m\}$ and $Y = \{y_1, \ldots, y_m\}$

Find

- a permutation matrix $\pi \in \Pi_m$,

$$\Pi_m := \left\{ \pi | \pi \in \{0, 1\}^{m \times m} \text{ where } \pi 1_m = 1_m, \pi^\top 1_m = 1_m \right\}$$

such that $\left\{ (x_i, y_{\pi(i)}) \text{ for } 1 \leq i \leq m \right\}$ is maximally dependent.

## Hilbert-Schmidt Independence Criterion

- HSIC is the square of the Hilbert-Schmidt norm of the cross covariance operator

$$\mathcal{C}_{xy} = \mathbf{E}_{xy}[(\phi(x) - \mu_x) \otimes (\psi(y) - \mu_y)],$$

where $\mu_x = \mathbf{E}[\phi(x)], \mu_y = \mathbf{E}[\psi(y)]$.

- In term of kernels, HSIC can be expressed as

$$\text{HSIC} = \|\mathcal{C}_{xy}\|_{\text{HS}}^2 = \mathbf{E}_{xx'yy'}[k(x, x')l(y, y')] + \mathbf{E}_{xx'}[k(x, x')]\mathbf{E}_{yy'}[l(y, y')] - 2\mathbf{E}_{xy}[\mathbf{E}_{x'}[k(x, x')]\mathbf{E}_{y'}[l(y, y')]].$$

- A biased estimator of HSIC given finite sample $Z = \{(x_i, y_i)\}_{i=1}^m$ drawn from $\text{Pr}_{xy}$ is

$$\widehat{\text{HSIC}} = (m-1)^{-2} \text{tr} HKHL = (m-1)^{-2} \text{tr} \tilde{K}\tilde{L},$$

where $K \in \mathbb{R}^{m \times m}, K_{ij} = k(x_i, x_j)$
$\tilde{K} := HKH.$

Advantages of HSIC are

- Computing HSIC is simple: only the kernel matrices $K$ and $L$ are needed;
- HSIC satisfies concentration of measure conditions, i.e. for random draws of observation from $\text{Pr}_{xy}$, HSIC provides values which are very similar;
- Incorporating prior knowledge into the dependence estimation can be done via kernels.

## Kernelized Sorting

Optimization problem

$$\pi^* = \text{argmax}_{\pi \in \Pi_m} \left[ \text{tr} K\pi^\top L\pi \right]$$

Sorting as a special case
For scalar $x_i$ and $y_i$ and a linear kernel on both sets, we can rewrite the optimization problem

$$\pi^* = \text{argmax}_{\pi \in \Pi_m} \left[ X^\top \pi Y \right]^2$$

This is maximized by sorting $X$ and $Y$.

## Optimization

Convex Objective and Convex Domain
- Define $\pi$ as a doubly stochastic matrix,

$$P_m := \left\{ \pi \in \mathbb{R}^{m \times m} \text{ where } \pi_{ij} \geq 0 \text{ and } \sum_i \pi_{ij} = 1 \text{ and } \sum_j \pi_{ij} = 1 \right\}$$

- The objective function $\text{tr} K\pi^\top L\pi$ is convex in $\pi$.

Convex-Concave Procedure
- Compute successive linear lower bounds and maximize

$$\pi_{i+1} \leftarrow \text{argmax}_{\pi \in P_m} \left[ \text{tr} K\pi^\top L\pi_i \right]$$

This will converge to a local maximum.
- Initialization is done via sorted principal eigenvector.

## Related Work

Instead of HSIC, we can use Mutual Information to measure the dependence between random variables $x_i$ and $y_{\pi(i)}$. MI is defined as, $I(X, Y) = h(X) + h(Y) - h(X, Y)$. We can approximate MI maximization by maximizing its lower bound. This then corresponds to minimizing an upper bound on the joint entropy $h(X, Y)$.

Optimization problem

$$\pi^* = \text{argmin}_{\pi \in \Pi_m} \left[ \log |HJ(\pi)H| \right],$$

where $J_{ij} = K_{ij}L_{\pi(i),\pi(j)}$.

This is related to the optimization criterion proposed by Jebara (2004) in the context of aligning bags of observations by sorting via minimum volume PCA.

## Applications

### Image Layout



Layout of 284 images into a 'NIPS 2008' letter grid using Kernelized Sorting. Gaussian RBF kernel is used for the image objects and also for the positions of the grid. One can see that images are laid out in the letter grid according to their color grading.



Layout of 320 images into a $16 \times 20$ grid using Kernelized Sorting.
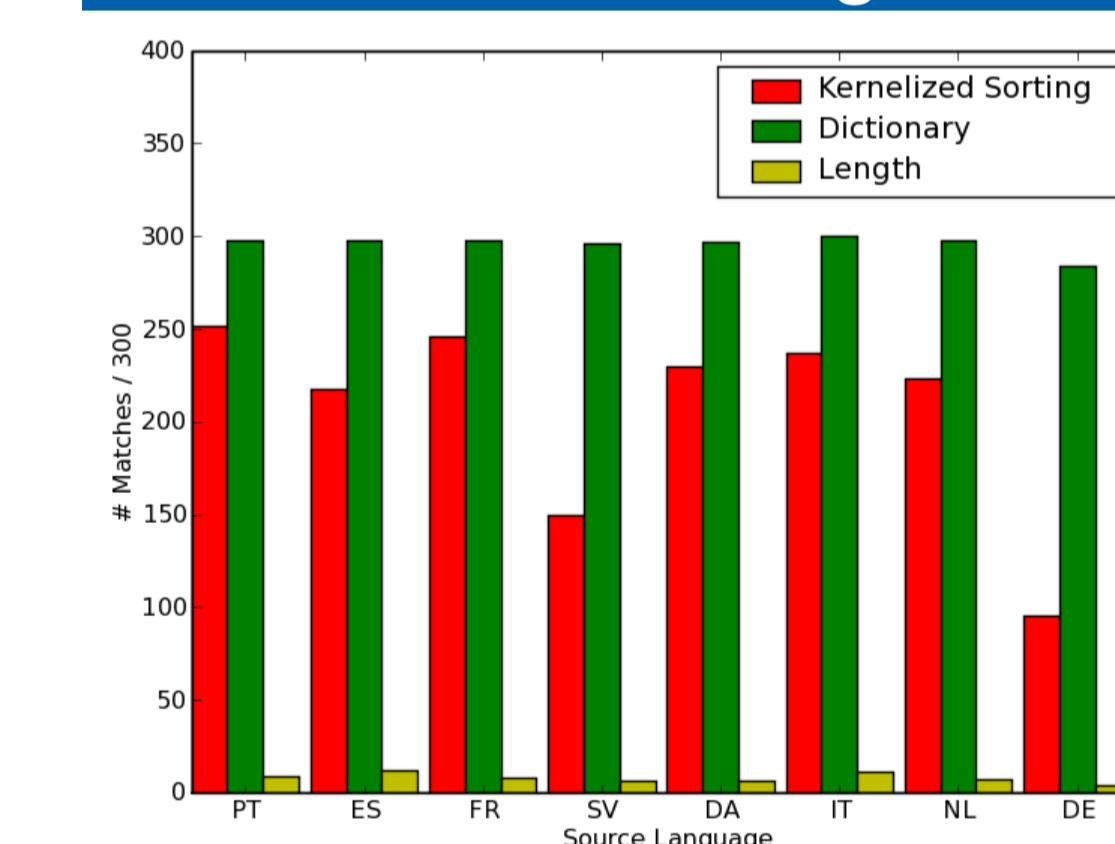


Layout of 320 images into a $16 \times 20$ grid using Generative Topographic Mapping. A compressed representation of images is shown. GTM does not guarantee unique assignments of images to nodes.
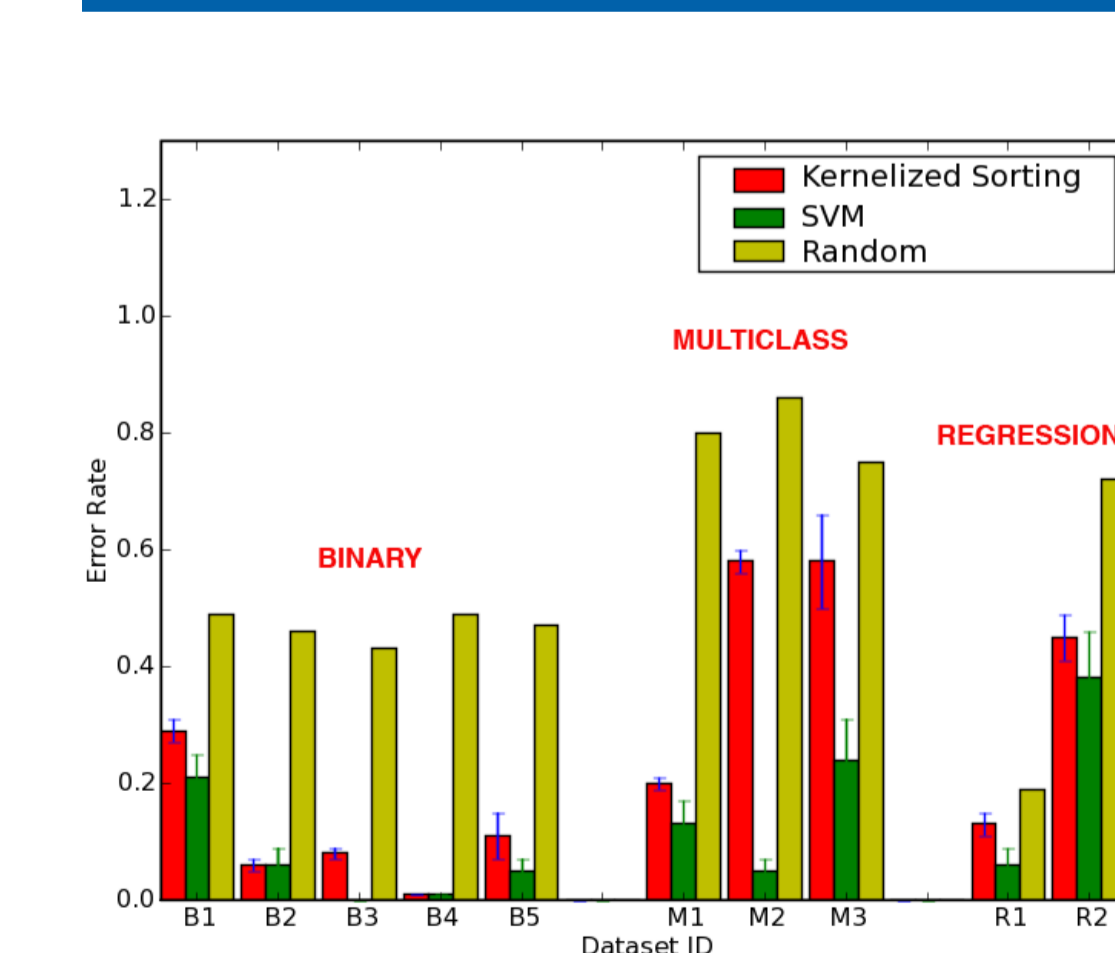
### Image Matching



Image matching as obtained by Kernelized Sorting. The images are cut vertically into two equal halves and Kernelized Sorting is used to pair up image halves that originate form the same images. 140 pairs out of 320 are correctly matched. This is quite respectable given that chance level would be 1 correct pair.

### Multilingual Document Matching



- Matching non-English documents (source languages) to its English translation (target language) of the Europarl Parallel Corpus.
- Standard TF-IDF features of a-bag-of-words kernel is used for Kernelized Sorting.
- Result comparisons with length based and dictionary based document matching.

### Data Attribute Matching



- Matching data attributes (or dimensions) of binary, multiclass, and regression data sets from UCI repository and LibSVM site.
- Homogenous misclassification loss and squared loss are used to measure performance.
- Gaussian RBF kernel is used for Kernelized Sorting.
- Results comparisons with random permutations and SVM classification/regression on original un-splitted data sets.

## Summary

- We generalize sorting by maximizing dependency between matched pairs of observations via HSIC.
- Applications of our proposed sorting algorithm range from data visualization to image, data attribute and multilingual document matching.