

Approximated Neighbours MinHash Graph Node Kernel

Nicolò Navarin and Alessandro Sperduti *

Department of Mathematics, University of Padova
via Trieste 63, Padova, Italy

Abstract. In this paper, we propose a scalable kernel for nodes in a (huge) graph. In contrast with other state-of-the-art kernels that scale more than quadratically in the number of nodes, our approach scales linearly in the average out-degree and quadratically in the number of nodes (for the Gram matrix computation). The kernel presented in this paper considers neighbours as sets, thus it ignores edge weights. Nevertheless, experimental results on real-world datasets show promising results.

1 Introduction

In the last few years, there has been a growing interest in machine learning techniques for graph-structured domains. For example, recommendation systems can be seen as link-prediction systems in bipartite graphs where nodes represent users and items [1]. Moreover, social networks can be represented as a graph where nodes are the users and links are their interactions; in this setting, it is interesting to predict (i.e. to suggest to users) possible novel connections [2]. In these and other scenarios, the learning problem can be expressed as classification (or ranking) over the nodes in a single, huge graph. Kernel methods are popular learning algorithms in the machine learning community, thanks to their state-of-the-art performances in many real-world tasks and to the strong theoretical guarantees behind the selected solution given by statistical learning theory. However, in many applications, their diffusion suffers from the computational complexity that, while not being high in general, may nonetheless be too high for the data amount available nowadays. When data can be represented in a structured form (in this paper we deal with nodes in a graph), the kernels available in literature scale more than quadratically in the number of nodes. Note that we consider graph node kernels, that are defined over the nodes in a single graph. This setting is different with respect to the one where each example is a different graph, e.g. the one considered in [3], and thus it requires different techniques. In this paper, we propose a novel kernel function for graph nodes that: (i) is efficient to compute; (ii) can be computed on demand on a pair of nodes at a time; (iii) performs well on different and heterogeneous real-world graph datasets. The idea is to define the kernel on (summaries of) the local neighbourhood of two nodes in the graph.

*This work was supported by the University of Padova under the strategic project BIOINFOGEN.

2 Definitions and notation

A kernel method is composed by the following two modules: (i) a learning algorithm, that expresses the solution only via dot products between training examples; (ii) a kernel function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, i.e., a symmetric positive semi-definite function that corresponds to a dot product in a Reproducing Kernel Hilbert Space (RKHS), which means that there exists a $\phi : \mathcal{X} \rightarrow \mathcal{K} \subseteq \mathbb{R}^D$, where \mathcal{K} is an Hilbert space (commonly referred to as *feature space*), such that $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ with $x_i, x_j \in \mathcal{X}$. Note that: (i) the input space \mathcal{X} can be any space; (ii) we do not need to know ϕ explicitly.

Given a training set $Tr = \{x_1, \dots, x_n\}$, the Gram matrix is a matrix with entries $\mathbf{K}_{ij} = k(x_i, x_j)$.

Let $G = (V, E)$ be a graph, where $V = \{v_1, \dots, v_n\}$ is the set of nodes ($|V| = n$) and $E = \{(v, u) | u, v \in V\}$ is the set of edges ($|E| = m$). Every edge can have an associated weight $w_{u,v}$, that is zero by convention if $(u, v) \notin E$. Note that, in our setting, \mathcal{X} is the space of nodes in the considered graph, and the labels we aim to predict are associated to the nodes in the graph. Let $\rho = \frac{1}{n} \sum_{u \in V} |\{(u, v) | (u, v) \in E\}|$ be the average out-degree of nodes in G . Given a graph G , let us define its adjacency matrix \mathbf{A}_G with entries $\mathbf{A}_{G_{ij}} = w_{ij}$. For the sake of notation, in the following we will omit the subscript G . Let us define the Laplacian matrix as $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where \mathbf{D} is a diagonal matrix with entries set to the summation over the corresponding row of the adjacency matrix, i.e. $\mathbf{D}_{ii} = \sum_j \mathbf{A}_{ij}$.

3 Related work

In the following, we will give the definition of four state-of-the-art functions that, given a graph G , compute the Gram matrix of the nodes in G . We will consider these functions as baselines in the experimental section.

All the kernels are based on random walks. These kernels are relatively fast to compute, and show good predictive performance [1, 4].

Laplacian Exponential Diffusion kernel (LEDK) [5] is defined as $LEDK(G) = e^{-\beta \mathbf{L}}$, where $0 \leq \beta \leq 1$ is a parameter used to control the rate of diffusion.

Markov Diffusion kernel (MDK) [1], depends on a parameter t controlling the length of the considered walks. Let us define the probability transition matrix \mathbf{P} as $\{\mathbf{P}\}_{ij} = \mathbf{A}_{ij}/\mathbf{D}_{ii}$. Moreover, let us define $\mathbf{Z}(t) = \frac{1}{t} \sum_{\tau=1}^t \mathbf{P}^\tau$. Then, we can define the kernel as: $MDK(G) = \mathbf{Z}(t)\mathbf{Z}^T(t)$.

Markov Exponential Diffusion kernel (MEDK) [6]. Let us define $\mathbf{M} = (\mathbf{D} - \mathbf{A} - n\mathbf{I})/n$ where n is the number of nodes in the graph. Then, MEDK is defined as $MEDK(G) = e^{-\beta \mathbf{M}}$.

Regularized Laplacian Kernel (RLK) [5]. Is defined as $RLK(G) = \sum_{n=1}^{\infty} \beta^n (-\mathbf{L})^n = (\mathbf{I} + \alpha \mathbf{L})^{-1}$ where $\alpha > 0$ is a parameter.

The kernels LEDK and MEDK require to compute a *matrix exponential*, that (e.g. using Sylvester's equations) involves $O(n^3)$ operations. The kernel MDK computes t matrix multiplications, each one with a cost of approximately

$O(n^3)$ (the fastest algorithm runs in $O(n^{2.373})$). The kernel RLK requires a matrix inversion (same complexity as matrix multiplication). The computational complexities of the kernels presented in this section may be prohibitive when dealing with real-world graphs containing hundreds of thousands, or million of nodes. In this paper, we propose a kernel for graph nodes based on MinHash. The computation of the Gram matrix scales only quadratically in the number of nodes and linearly in the average node out-degree¹ ρ . MinHash has been previously applied to graph kernels, e.g. in [7] but, to the best of our knowledge, it has never been applied to graph node kernels.

3.1 MinHash

Let A and B be two sets, with elements in a domain Σ . The Jaccard similarity, a commonly adopted similarity measure between sets, is defined as the ratio between the number of elements in the intersection and the number of elements in the union, i.e. $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$. MinHash [8] is an algorithm, originally designed for near-duplicate detection in documents, that computes an unbiased estimation of the Jaccard similarity between sets of arbitrary sizes in linear time using a small and fixed memory space². Let $h : \Sigma \rightarrow N_D$ be a perfect hash function³ that maps the elements of Σ to the set $N_D = \{x | x \in \mathbb{N}, x \leq 2^D - 1\}$ of *distinct* integers represented by at most D bits. Let us define, for a set S , $h_{min}(S)$ to be the minimal member of S with respect to h . If we apply h_{min} to both A and B , we get the same value exactly when the element of the union $A \cup B$ with minimum hash value lies in the intersection $A \cap B$. The probability of this being true is: $Pr[h_{min}(A) = h_{min}(B)] = \frac{|A \cap B|}{|A \cup B|} = J(A, B)$. However, this estimator has a too high variance to be useful. The idea of the MinHash scheme is to reduce this variance by averaging together several variables constructed in the same way. Moreover, we can consider hash functions with collisions since they generally have only a modest influence on our estimate (and the estimation remains still unbiased). Let k be a fixed integer (a parameter), and let us define a Minwise independent family of hash functions, $\{h^1 \dots h^k\}$. The MinHash algorithm computes the following value: $MH(A, B) = \frac{1}{k} \sum_{i=1}^k \delta(h_{min}^i(A), h_{min}^i(B))$, where δ is the Kronecker delta function. MH is an average of unbiased estimators, thus it is in turn an unbiased estimator of $J(A, B)$. The expected error is, using Chernoff bounds, $O(1/\sqrt{k})$.

4 Methods

Let $sp(u, v)$ be a function that returns the length of a shortest path between nodes u and v . We can define $N^r(v) = \{u | sp(v, u) = r\}$ to be the neighbourhood of radius r (more compactly r -neighbours) of the node v in a graph G . Note

¹In the worst case, ρ is bounded by n ; however, in practical applications, it is much smaller.

²We adopted the MinHash implementation available at <https://github.com/ekzhu/datasketch>.

³A hash function with no collisions.

that $N^0(v) = v$ and $N^1(v) = \{u | (v, u) \in E\}$. We can define a kernel for the nodes in a graph as: $K(u, v) = \sum_{i=0}^r J(N^i(u), N^i(v))$, where $r \in \mathbb{N}$ is a fixed constant (hyper-parameter). Note that computing $J(A, B)$ between two sets A and B is linear in $\min(|A|, |B|)$. The Jaccard similarity is exactly the Tanimoto kernel between two sets. The proof of positive-semidefiniteness of the Tanimoto kernel can be found in [9].

Note that for computing the r -neighbours for all the nodes in a graph we have to perform n breadth-first visits of depth r , that are really slow when the graph is large. For this reason, we propose an approximation of the kernel just defined, using a scheme inspired from the Weisfeiler-Lehman method [10].

We start computing $N^1(v)$ for each $v \in V$. The complexity of this step is $O(n\rho)$. We then compute an approximation of $N^2(v)$ as $\hat{N}^2(v) = \bigcup_{u \in N^1(v)} N^1(u)$. We can easily extend this definition to arbitrarily large r as $\hat{N}^r(v) = \bigcup_{u \in \hat{N}^{r-1}(v)} N^1(u)$. This procedure is more efficient with respect to the computation of the r -neighbours because we have only to query the graph once for node neighbours. We can store each node's neighbours as a set and store them in a dictionary for efficiency purposes. The other operations are just set unions, that are generally fast. The complexity of computing the $\hat{N}^r(v)$ for each $v \in V$ becomes $O(n\rho^r)$. We can define the resulting kernel as: $K(u, v) = \sum_{i=0}^r J(\hat{N}^i(u), \hat{N}^i(v))$.

We can now introduce another approximation in the computation of J , using the MinHash technique presented in Section 3.1:

$$K(u, v) = \sum_{i=0}^r MH(\hat{N}^i(u), \hat{N}^i(v)). \quad (1)$$

Note that, differently to the other consider kernels, our proposal does not exploit edge weights.

MinHash allows to compute an unbiased estimation of the Jaccard similarity between two sets in a computational time that does not depend on the set cardinality. Instead, it depends just on the number of hash functions that are used.

Theorem 1. *The kernel in eq. 1 is positive semidefinite.*

Proof. If we consider the relation mapping a node u to the set of its $\hat{N}^i(u)$ according to the described procedure, the kernel in eq. 1 is an instance of the Convolution Kernels framework [11]. A previous result ([12], Theorem 2) proves that MH is a valid kernel over sets. \square

The strongest point of the proposed kernel is that its computation between two nodes is independent from the others. This property is particularly useful when the size of the dataset is big, such that storing the Gram matrix ($O(n^2)$ space) is unfeasible. In this case, it is possible to just compute the $\hat{N}^i(u)$ for each node in the graph. The cost of this step, for the whole graph, is $O(n\rho^r)$. Note that memory occupation of this data structure is linear with respect to the number of nodes and the considered neighbourhood (the r parameter). Then, one can compute the kernel between two nodes on-demand in a time that is

Dataset/kernel	LEDK	MDK	RLK	ANMHK
IMDB	73.05 ± 2.39	78.19 ± 0.44	66.58 ± 1.50	77.85 ± 0.54
IMDB_ALL	58.52 ± 1.17	72.24 ± 0.45	72.36 ± 2.17	79.44 ± 0.23
INDUSTRY_PR	34.28 ± 0.41	36.72 ± 0.34	34.53 ± 0.40	34.94 ± 0.19
INDUSTRY_YH	31.99 ± 0.38	46.89 ± 0.42	30.51 ± 0.41	37.35 ± 0.50
CORA	63.88 ± 0.67	75.85 ± 0.12	70.64 ± 0.24	72.22 ± 0.13
WEBKB_CORNELL	42.86 ± 0.42	53.99 ± 1.02	41.31 ± 0.02	56.76 ± 0.87
WEBKB_WISCONSIN	53.42 ± 0.95	70.34 ± 1.10	43.78 ± 0.09	74.52 ± 0.54
WEBKB_WASHINGTON	46.07 ± 1.36	65.87 ± 0.70	42.51 ± 1.74	68.50 ± 0.26
WEBKB_TEXAS	58.08 ± 0.97	68.11 ± 0.82	48.22 ± 0.01	73.10 ± 0.35

Table 1: Experimental results comparing the proposed kernel (ANMHK) vs other state-of-the-art ones on 9 real-world datasets.

linear in r and k . This approach is not possible with the other considered kernels: indeed, they can only compute the Gram matrix, via transformations of the adjacency matrix. If r and ρ are small with respect to n , we can consider them constants. Moreover, k is a constant because it is fixed a-priori. The complexity of the kernel computation between two nodes, under these assumptions, is constant, and the complexity of computing the whole Gram matrix is $O(n^2)$.

5 Experimental assessment

In the following, we describe the adopted datasets, originally used in [13]⁴.

IMDB (IMDB_ALL): dataset of movies released in the US, with class labels identifying whether the opening weekend box-office receipts will exceed \$2M. Two movies share a link if they share a production company (a production company, producer, director, or actor).

CORA: this dataset comprises computer science research papers. It includes the full citation graph as well as labels for the topic of each paper. There are seven possible labels.

INDUSTRY_YH (INDUSTRY_PR): this dataset contains companies that are linked via co-occurrence in business stories from the web (Newswire press releases). An edge was created between two companies if they appeared together in a document. The resulting network comprises 1798 (2189) companies. The labels of the companies are based on Yahoo!'s 12 industry sectors.

WEBKB⁵: this dataset consists of sets of web pages from four computer science departments (Cornell, Wisconsin, Washington, Texas), where each page has been manually labeled into 7 categories: course, department, faculty, project, staff, student, or other. The edges are based on co-citation.

Since the MEDK kernel performed poorly on all the considered datasets, for the sake of space, we decided not to report its results. In Table 1, the accuracy results of the proposed kernel and the other three state-of-the-art kernels on the above described datasets are reported. The results are obtained with a 10-fold cross validation (CV) that uses just 1 fold for training and hyper-parameter optimisation (with a nested 5-fold CV), and 9 for testing, i.e. 90%

⁴We retrieved the datasets from <http://netkit-srl.sourceforge.net/data.html>

⁵<http://www.cs.cmu.edu/~webkb/>

of the data labels have to be predicted. As a learning algorithm, we adopted a one-versus-one multiclass SVM with $C \in \{10^{-4}, \dots, 10^4\}$. The kernel parameters have been optimised in: $\beta \in \{0.01, \dots, 0.05, 0.1, 0.3, 0.5\}$ for LEDK and MEDK; $\alpha \in \{10^{-2}, \dots, 10^3\}$ for RLK; $t \in \{1, 2, 3, 5, 10, 50, 100\}$ for MDK; $r \in \{1..10\}$ for the proposed ANMHK.

The proposed method gets the best performance on five datasets. On the other four datasets, it ranks in second place. These results suggest that kernels for graph nodes based on local node neighbourhood can be expressive enough for several real-world tasks.

6 Conclusions and future work

We proposed a novel kernel for graph nodes that scales linearly in the average node out-degree of the graph. The kernel presented in this paper considers neighbours as sets, thus it ignores edge weights. We plan to modify the kernel including weighted edges in a future extended version of this paper.

References

- [1] François Fouss, Kevin Francoisse, Luh Yen, Alain Pirotte, and Marco Saerens. An experimental investigation of kernels on graphs for collaborative recommendation and semi-supervised classification. *Neural Networks*, 31:53–72, 2012.
- [2] David Liben-nowell and Jon Kleinberg. The Link Prediction Problem for Social Networks. In *CIKM*, pages 556–559, 2003.
- [3] Giovanni Da San Martino, Nicolò Navarin, and Alessandro Sperduti. Ordered Decompositional DAG Kernels Enhancements. *Neurocomputing*, 192:92–103, 2016.
- [4] Risi Imre Kondor and John Lafferty. Diffusion kernels on graphs and other discrete structures. In *ICML*, 2002.
- [5] Alexander J Smola and Risi Kondor. Kernels and Regularization on Graphs. In *Learning theory and kernel machines.*, pages 144–158. Springer Berlin Heidelberg, 2003.
- [6] BoLin Chen, Min Li, JianXin Wang, and Fang-Xiang Wu. Disease gene identification by using graph kernels and Markov random fields. *Science China Life Sciences*, 57(11):1054–1063, 2014.
- [7] Carlos H.C. Teixeira, Silva Arlei, and Meira Jr. Wagner. Min-Hash Fingerprints for Graph Kernels: A Trade-off among Accuracy, Efficiency, and Compression. *Journal of Information and Data Management*, 3(3):227–242, 2012.
- [8] Andrei Z Broder. Identifying and filtering near-duplicate documents. In *Combinatorial Pattern Matching*, pages 1–10. Springer Berlin Heidelberg, 2000.
- [9] Liva Ralaivola, Sanjay J Swamidass, Hiroto Saigo, and Pierre Baldi. Graph kernels for chemical informatics. *Neural networks*, 18(8):1093–110, oct 2005.
- [10] Giovanni Da San Martino, Nicolò Navarin, and Alessandro Sperduti. Graph Kernels Exploiting Weisfeiler-Lehman Graph Isomorphism Test Extensions. In *Neural Information Processing*, pages 93–100. Springer, 2014.
- [11] David Haussler. Convolution Kernels on Discrete Structures. Technical report, Department of Computer Science, University of California at Santa Cruz, 1999.
- [12] Ping Li, Anshumali Shrivastava, Joshua Moore, and Arnd Christian Konig. Hashing Algorithms for Large-Scale Learning. *Advances in Neural Information Processing Systems 24*, pages 2672–2680, 2011.
- [13] Sofus A Macskassy and Foster Provost. Classification in Networked Data: A Toolkit and a Univariate Case Study. *JMLR*, 8(December 2004):935–983, 2007.