

Graph Kernels and Support Vector Machines for Pattern Recognition

Léo Andéol*

Master DAC, Sorbonne Université
Paris, France

Supervised by: Prof. Hichem Sahbi

May 2019

Abstract

The problem of having a framework to classify graphs is becoming increasingly important in the era of data. The issue has been tackled since the beginning of the millennium and there have been significant progress. This report introduces all necessary knowledge to understand the topic and then reviews different graph kernels while focusing on a random walk kernel and its optimization. Some experiments are then conducted to verify the information given by the state of the art, and some attempts are made to accelerate the different methods to compute the kernel.

*leo.andéol@gmail.com

Contents

1	Introduction	3
2	Methodology	4
2.1	Background	4
2.1.1	Graphs	4
2.1.2	Support Vector Machines	5
2.1.3	Kernels	7
2.2	State of the Art	8
2.2.1	Graph Kernels	8
2.2.2	Graphlets	8
2.2.3	Random Walks	8
2.2.4	Digest	14
3	Experiments	15
3.1	Databases and Metrics	15
3.1.1	A synthetic graph database	15
3.1.2	Metrics	16
3.2	Evaluation Protocols	17
3.3	Performance	19
3.3.1	Performance on synthetic data	20
3.3.2	Performance on real datasets	22
4	Conclusion and Future Work	23
A	Appendix	24

1 Introduction

Nowadays the world is shifting towards a new era where data is the most valuable resource, and researchers are working on algorithms to exploit it the best possible way. Most research is currently focused on the so-called "big data", huge quantities of data which are very sparse and of poor quality. However, there also exist structured databases of good quality, in the form of graphs, this very form being the one studied throughout this project. Graphs can be very useful in various fields, but are partly known for their use in biology, as they can be used to represent proteins or other types of molecules.

Classification of such graphs is an important problem of pattern recognition which affects a lot of sectors such as medicine, biology, and more. It was studied since the beginning of the millennium and significant progress has been made : several methods were discovered and gave good results but met a big issue, the complexity of these methods. Indeed, they do not scale well to large graphs as their complexity tends to be several times polynomial. Since then, the main objective of research on this topic has been to either find more computationally efficient algorithms, or to find methods to accelerate the ones already in use.

The first part of this report will introduce all the necessary background knowledge which are graphs as already mentioned, but also Support Vector Machines (SVMs) which is a powerful classification algorithm, and kernels which are transformations of data usually used to increase the accuracy of SVMs. However, graphs being non-vector data, kernels have also the purpose of becoming a metric of comparison between two instances of graphs. Thus, the main principles of graph kernels will be introduced as well as their definitions. Afterwards, the different methods used to improve the computation complexity of the kernel will be introduced together with their own complexities, advantages and disadvantages.

The second part of this report will be focused on the technical side and experiments conducted during the project. A synthetic graph database made of toy data was required to conduct simple and quick experiments. Indeed, it was made in order to verify the claims made in the main publication studied thanks to easy control on matters such as labeled and unlabeled graphs, and variations of the size of either the database or the graphs inside of it. Different challenges and problems met during the implementation will be explained, as well as their solution, if one has been found. Then, the main subject of this part will be discussed : the accuracy and computation time of different methods, estimation of their complexities and different attempts to approximate those methods. Finally, an experiment on real datasets of enzymes and proteins is conducted, and the results analyzed.

introduction
générale plus
étendue : 1
et 1/2 page
: donner un
aperçu/résumé
du rapport
: un para-
graphe par
élément -
ECRIRE A
LA FIN

garder som-
maire ? en
général il n'y
en a pas

2 Methodology

2.1 Background

2.1.1 Graphs

Revoir
toutes les
refs

Definition A graph[1] is a type of mathematical structure that represents connections between objects. It is more precisely an ordered pair $G = (V, E)$ of two sets: vertices V (or nodes) and edges E that connect two vertices together.

$$E \subseteq \{(u, v) : (u, v) \in V^2\} \quad (1)$$

Vertices represent objects and are usually depicted as circles or spheres whereas edges link pairs of vertices.

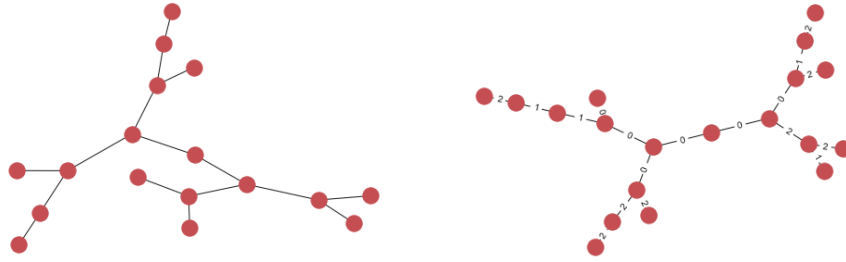


Figure 1: Two "tree" graphs, resp. unlabeled and labeled

Properties Graphs have a lot of properties that will be used in the rest of this project and will be introduced here

- A graph is undirected if and only if : $\forall (u, v) \in E \implies (v, u) \in E$
- $|V|$ is the number of vertices and $|E|$ the number of edges
- The degree of a vertex $d(v)$ is the number of vertices it is connected to, or the number of distinct edges connected to it.
- A path is a sequence of connected edges linking distinct vertices.
- A cycle is a path where the first and last vertex is the same.
- A graph can have labels (sometimes called colors) either on its vertices or edges (or both). They can take various forms such as being integers, more generally elements of a finite or infinite set and even continuous (such as $\in \mathbb{R}$).

- A graph is said to be connected if all vertices are connected by a path with any other vertex.
- A graph is said to be a tree if it is connected and doesn't have any cycle.
- A line graph $G' = (V', E')$ of a graph $G = (V, E)$ is composed of a set of vertices V' of a new vertex for each unique $e \in E$ and a set of edges E' where there is an edge $(u, v) \in V'^2$ if and only if there was a vertex connecting the two former edges u and v together.
- A subgraph $G' = (V', E')$ of a graph $G = (V, E)$ is that graph restricted to a subset of nodes $V' \subseteq V$ while keeping all the edges with those vertices $E' = \{(u, v) : (u, v) \in E \wedge (u, v) \in V'^2\} \subseteq E$

These methods are focused on undirected unlabeled graphs, and undirected graphs labeled by finite sets. Moreover, only edge-labeled graphs will be studied, thus if a graph is node-labeled, simply taking its line graph will do.

Adjacency Matrix The adjacency matrix of an undirected graph represents the presence or absence of an edge between two specific vertices. It is defined as the matrix A of dimension $|V| \times |V|$ where entries are given by

$$A_{i,j} = \begin{cases} 1 & \text{if } i \neq j \text{ and } (i, j) \in E \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

History Graphs were first used in their modern form to represent the problem of Seven Bridges of Königsberg, and have been ever since used to represent maps, and thus path-finding algorithms were developed. They have also been used to represent flow problems, scheduling problems, networking routing and many others. Graphs can also be used, with labels, to represent different types of molecules and interactions between them, or more simply to represent molecules by considering atoms as vertices and bonds as edges.

The next part will introduce Support Vector Machines, the algorithm used in this project to classify new unknown graphs into known categories. It has been shown[2] that it can be used effectively on proteins and enzymes.

2.1.2 Support Vector Machines

Support Vector Machines (SVMs)[3] are a type of machine learning algorithms discovered in the early 90s[4]. It was originally a classification algorithm however it has been expanded since to regression and clustering too.

Classification Classification is one of the two problems of supervised learning that aims to automatically recognize and classify observations \mathbf{x} , such as recognizing handwritten digits. It can be split into two tasks where the first task is to learn from training data how to decide in which class (or group) to classify an example from its attributes and the second is to predict the class

of new examples. All of that, while minimizing errors. This problem can be formalized as

Definition 1. Classification is the problem of finding the best function $f : \mathbb{R}^d \rightarrow \{0..k\}$ among a set of functions F while minimizing a risk function R

$$f^\star = \operatorname{argmin}_f R(f) \quad (3)$$

where k is the number of classes of the problem and d the number of features of data

SVMs are especially powerful at this task and widely used for the two following reasons.

Margin and support vectors The Support Vector Machines are based on the model of the Perceptron[5], another classification algorithm that tried to find an hyperplane that discriminates the two sets. The main issue with the Perceptron is that it has no formal guarantee to find an optimal hyperplane, and usually would not find it. Moreover, if the data were not linearly separable, the algorithm would not converge.

In order to tackle these issues, the SVM offered three fixes. First, a margin was added in the loss L function in order to not only properly classify the data samples, but also with the maximum certainty, i.e. as far as possible from the decision boundary.

$$L(y_i, \mathbf{x}_i, w) = \max(0, -y_i \mathbf{x}_i \cdot \mathbf{w} + w_0) \implies L(y_i, \mathbf{x}_i, w) = \max(0, 1 - y_i \mathbf{x}_i \cdot \mathbf{w} + w_0) \quad (4)$$

Where \mathbf{x} is the feature vector of an instance, y its class and \mathbf{w} and w_0 respectively the weight vector and bias calculated by the algorithm. Then, seeing this solution was not enough, as there was still an infinity of possible solutions, it was proven that the margin γ between points and the decision boundary was inversely proportional to the norm of the weight vector $\|\mathbf{w}\|$.

$$\gamma_i = \frac{y_i(\mathbf{x}_i \cdot \mathbf{w} + w_0)}{\|\mathbf{w}\|} \quad (5)$$

These two fixes gave us the first version the SVM, the hard-margin SVM

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + w_0) \geq 1 \quad \forall i \in \{1..n\} \quad (6)$$

However, if the data were not linearly separable the algorithm would not be optimizable since the quadratic programming problem requires all points to be correctly classified. Then, a new term ξ was introduced as an error tolerance, as well as a factor C that would determine the balance between error tolerance and minimization of the weight vector's norm $\|\mathbf{w}\|$.

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall i \in \{1..n\} \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + w_0) \geq 1 - \xi_i \end{aligned} \quad (7)$$

Then, the next section will introduce the second advantage of SVMs.

Parler de la
theorie de
l'apprentissage
? risque
structurel
etc ? et dire
quelle loss
on utilise ou
pas ?

2.1.3 Kernels

In its dual form, the SVM problem only requires a dot product between the observations' vectors.

$$\begin{aligned} \max \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \forall i \quad 0 \leq \alpha_i \leq C \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned} \quad (8)$$

Where n is the number of instances. Then, a nonlinear transformation $\phi(x)$ can be used to replace the dot-product $\mathbf{x}_i \cdot \mathbf{x}_j$ by $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$, effectively augmenting the size of vectors from the observation space into a larger feature space and thus the expressiveness of the model by enhancing the separability of the data. However, for complex transformations, such as in infinite dimensions, the function ϕ is not defined. But it was found that the standard dot product can be replaced by another function as long as it is a bilinear positive semi definite form.

Definition 2. A kernel K is positive semi definite if and only if

$$\sum_{i=1}^n \sum_{j=1}^n \kappa(\mathbf{x}_i, \mathbf{x}_j) c_i c_j \geq 0 \quad \forall i \in \{1..n\} \quad c_i \in \mathbb{R} \quad (9)$$

Thus, the dot product can be replaced by a function $K(x_1, x_2)$ which will indirectly map the data to higher dimensions like the polynomial kernel¹⁰, even infinite with the RBF kernel¹¹ and return the dot product of those transformations, while remaining computable in all cases. This replacement is usually called the Kernel Trick. For example, the polynomial kernel is defined as follows

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + c)^d \quad (10)$$

where d is the degree of the polynomial and c the constant bias. Another example, the RBF kernel is defined as follows

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}} \quad (11)$$

where σ is an hyperparameter. Ever since the foundation of SVMs, the kernel trick became a big focus of the machine learning community as it naturally fits in the algorithm and allows supervised learning on very complex data, and enjoying greater accuracy than most algorithms.

The fact that the dot product can be replaced by a function without specifying the map ϕ is the reason why the Kernel Trick became a big focus of research. In particular, it will be possible to make a kernel for non vector data such as graphs, as shown in the next section.

Est-ce qu'il faut ren-
trer dans le
détail ? ra-
jouter quoi?

2.2 State of the Art

This section will review the state of the art kernels for graphs. Particular attention will be given to the random walk kernel.

2.2.1 Graph Kernels

Graph Kernels are a type of R-convolution kernels[6] applied to graphs, which are kernels that are based on decompositions of complex objects and comparisons of those decompositions.

Graph Kernels have been studied since Support Vector Machines started getting popular[7]. Since then, a lot of progress has been made, and several types of kernels have been discovered, such as random walk and graphlet[8] kernels, the two main families of kernels (and the ones studied in this section), but not the only ones. There are also graph kernels that use subtrees[9], shortest paths[10] and several others that are unfortunately usually either too complex to compute or not positive semi-definite[11]. However it has been shown[12] that in some cases non positive semi definite kernels can still be used efficiently either by using them as a dissimilarity measure. They can also be artificially made positive semi definite.

2.2.2 Graphlets

A graphlet is a non-empty graph of $k = 3, 4$ or 5 nodes. The idea of the graphlet kernel κ is to generate a set S all possible graphlets of a certain size k and calculate for each of them the frequency of occurrence in a certain graph $G = (V, E)$ to build a vector \mathbf{f}_G of size $|S|$. The kernel is then defined as follows

Definition 3. Let G and G_2 be two graphs, \mathbf{f}_G and \mathbf{f}_{G_2} the frequency vectors of respectively G and G_2 , then the kernel κ is defined as

$$\kappa(G, G_2) = \mathbf{f}_G^\top \mathbf{f}_{G_2} \quad (12)$$

However the complexity of computing this vector \mathbf{f}_G is $O(n^k)$ where $n = |V|$ and is thus computationally extremely expensive. In order to address this issue, a sampling method was introduced. Indeed, if enough samples are drawn, the Law of Large numbers indicates that the sampled distribution will converge to the actual (theoretical) distribution. However a bound has been found on the complexity of such sampling[13] allowing better precision.

2.2.3 Random Walks

Graph kernels based on Random Walks have been studied since kernels started to be used on SVMs[7], several of them were discovered and then united[14]. Random walks are a type of rather intuitive algorithms : the idea is to randomly walk through a graph and then compare it to random walks in another graph. Actually, all the random walks are computed at the same time by making the adjacency matrix stochastic by dividing each column of the matrix by

developper
plus les
autres ker-
nels?

comment les
méthodes
operent :
differentes
sous meth-
odes (ex
methode
avaec labels,
methode
avec ori-
enté ou pas)
 \Rightarrow partie
déjà assez
longue ?

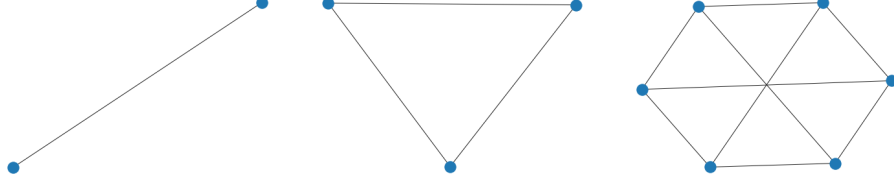


Figure 2: A graph G_1 , a graph G_2 and the product graph $G_1 \otimes G_2$

its sum (normalizing it) and using it as a Markov chain.

Moreover, it has also been shown[15] that performing a walk on two separate graphs at the same time is the same as performing a walk on the product graph. The product graph is in simpler terms a subset of the Cartesian product of the two graphs, taking all possible combinations of nodes and linking those nodes when there is an edge between two specific nodes in both the original graphs.

Definition 4. Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two graphs, the product graph $G_{\times} = (V_{\times}, E_{\times})$ is then defined as

- $V_{\times} = \{(v, w) : v \in V_1, w \in V_2\}$
- $E_{\times} = \{((v_1, w_1), (v_2, w_2)) : v_1, v_2 \in V_1^2, w_1, w_2 \in V_2^2\}$

For example, on very simple graphs, the following result is obtained

However, technically, the adjacency matrix W_{\times} of a such graph is the result of the Kronecker (tensor) product of the two adjacency matrices A and A_2 of the two graphs [16]

$$W_{\times} = A \otimes A_2 \quad (13)$$

In case the graphs are labeled, then the adjacency matrix W_{\times} is the sum of the Kronecker products of adjacency matrices restricted to each label l

$$W_{\times} = \sum_{l=1}^d A_1^{(l)} \otimes A_2^{(l)} \quad (14)$$

Where d is the size of the label set and $A_x^{(l)}$ the adjacency matrix of G_x restricted to the label li which is defined as follows

$$A^{(l)} = \begin{cases} 1 & \text{if } i \neq j \text{ and } (i, j) \in E \text{ and } (i, j) \text{ is labeled } l \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

The *vec* operator will also be used as it is very linked to the tensor product. It is the vector obtained by stacking all columns of A .

Kernel definition The idea of the kernel is to perform random walks simultaneously on two graphs in order to compare the walks they have in common. As shown above, that can be achieved by simply performing a random walk on the product graph. By making the adjacency matrix stochastic, and computing the k -th power of the matrix W_{\times} , the probabilities obtained are the ones of walks in common (starting and ending at the two same nodes) between the two graphs of length k . By taking into account the start p_{\times} and end q_{\times} probabilities, the sum of walks of all possible lengths can be calculated. Moreover using a power series $\mu(k)$ helps de-emphasizing longer walks and makes the sum converge.

Definition 5. Let p_{\times} and q_{\times} be respectively the start and end probability of each node, and let W_{\times} be the adjacency matrix of the product graph of G and G' , and finally $\mu(k)$ be a convergent function of k .

$$\kappa(G, G') = \sum_{k=0}^{\infty} \mu(k) q_{\times}^{\top} W_{\times}^k p_{\times} \quad (16)$$

This kernel can be computed in $O(kn^6)$ where k is the number of iterations

Inverse Kernel There is a specific case of the kernel following provided a power series is used. Indeed, by replacing the power series as follows $\mu(k) = \lambda^k$ the following kernel is obtained

$$\kappa(G_1, G_2) = q_{\times}^{\top} (I - \lambda W_{\times})^{-1} p_{\times} \quad (17)$$

It can also be computed in $O(n^6)$ however in practice the hidden factor is much lower which makes it a better option for very small graphs and databases and the accuracy much better since there is not any tradeoff on the number of iterations. Moreover, it will inspire some of the following methods.

Sylvester Equation A Sylvester equation, also sometimes called Lyapunov equation is a matrix equation of the following shape

Definition 6. Let A , B , and C be matrices of compatible shapes, then the Sylvester equation is

$$AX + XB = C \quad (18)$$

And the discrete-time Sylvester Equation is

$$AXB + C = X \quad (19)$$

Which can be generalized as

$$\sum_{i=0}^d A_i X B_i = X \quad (20)$$

Only the discrete-time Sylvester and its generalization will be used in this method and will be referred to as the Sylvester equation. These equations are solved usually using Schur decompositions in $O(n^3)$ for the basic Sylvester equation and in unknown time for the generalized version[14]. This equation can be used to solve our kernel faster assuming a geometric function $\mu(k) = \lambda^k$ with λ an hyperparameter. It can be achieved by replacing the A and B matrices by our adjacency matrices A_1 and A_2 , adding λ , and C by our start probabilities p_\times while vectorizing the whole equation thus obtaining

$$\text{vec}(M) = \text{vec}(\lambda A_1 M A_2) + p_\times \quad (21)$$

From which the following can be obtained (where I is the identity matrix)

$$q_\times^\top \text{vec}(M) = q_\times^\top (I - \lambda W_\times)^{-1} p_\times \quad (22)$$

Thus, having calculated the inverse kernel by an alternative method. The advantage of this method is that it is very fast, but unfortunately limited to unlabeled graphs, unless there is an implementation for the generalized Sylvester equation. Moreover, compared to the others, this option is very simple and does not require specific parameters, but like many others, suffers when it encounters singular matrices.

Conjugate Gradient Method The Conjugate Gradient Method[17] is an optimization algorithm used to find approximate solutions of linear systems that have a positive semi definite matrix. As its name suggests, it is a gradient based iterative algorithm. The main idea is that normal gradient descent only takes into account the gradient at the current step of the algorithm, and depending on the step size may cancel progress made in the previous step. The

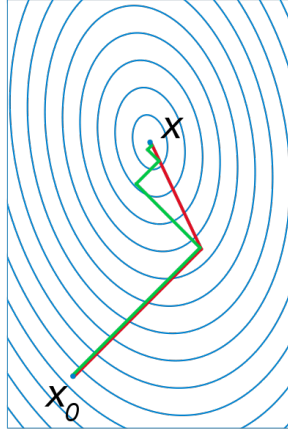


Figure 3: Conjugate gradient (red) compared to Gradient Descent (green)

Conjugate Gradient method was introduced as a fix for that issue. The idea of

the algorithm is to keep former gradients in order improve the convergence speed by making the new gradient orthogonal to all former ones through the Gram-Schmidt process, at each step. Thus, unless there are approximation errors, the algorithm guarantees convergence in n steps, being the number of dimensions of the problem, here the size of our graph.

However, this algorithm has limitations, it can only be used to solve linear systems of equations as $Ax = b$ under the constraint that A is symmetric positive definite. In our case, this algorithm can be used to solve the following problem

$$(I - \lambda W_{\times} x = p_{\times} \quad (23)$$

And the kernel value can be obtained by computing $q_{\times}^{\top} x$. The complexity of that method is $O(rn^3)$ for unlabeled graphs and $O(rdn^3)$ for labeled graphs where r is the effective rank of the matrix W_{\times} . It is important to note that this kernel also assumes a geometric $\mu(k)$ function.

This method is thus slower than the Sylvester Equation, however it is easily applied to labeled graphs, which is its main advantage.

Fixed Point Iterations A fixed point is a value for which a function, returns that same value

Definition 7. Let S a set and f a map $f : S \Rightarrow S$, $x \in S$ is a fixed point of f if and only if $x = f(x)$

Fixed point iterations is a method of computing a fixed point of a function by applying repeatedly the following equation until $\|x_{n+1} - x_n\| < \epsilon$ where ϵ is the acceptable level of error.

$$x_{n+1} = f(x_n) \quad (24)$$

There is also a guarantee of convergence

Theorem 1. *In order for the fixed-point iterations to converge, it is necessary that $\lambda < |\xi|^{-1}$ where ξ is the largest eigenvalue of W_{\times}*

le mettre en
appendix?

Proof. Let $x_0 = p_\times$ and $t \gg 0$

$$\begin{aligned}
x_{t+1} - x_t &= p_\times + \lambda W_\times x_t - x_t \\
&= p_\times + (\lambda W_\times - 1)(p_\times + \lambda W_\times x_{t-1}) \\
&= p_\times - p_\times + \lambda W_\times p_\times + (\lambda W_\times)^2 x_{t-1} - \lambda W_\times x_{t-1} \\
&= \lambda W_\times (p_\times - x_{t-1} + \lambda W_\times x_{t-1}) \\
&= \lambda W_\times (p_\times - \lambda W_\times x_{t-2} - p_\times + \lambda W_\times (\lambda W_\times x_{t-2} + p_\times)) \\
&= \lambda W_\times (\lambda W_\times (\lambda W_\times x_{t-2} + p_\times - x_{t-2})) \\
&= (\lambda W_\times)^2 (\lambda W_\times x_{t-2} + p_\times - x_{t-2})
\end{aligned}$$

A pattern is found

$$\implies x_{t+1} - x_t = (\lambda W_\times)^t (\lambda W_\times x_0 + p_\times - x_0)$$

And because $x_0 = p_\times$

$$= (\lambda W_\times)^{t+1} p_\times$$

Which decreases only when

$$\xi_1 < 1$$

The largest magnitude eigenvalue of λW_\times

$$\implies \lambda < |\xi_1|^{-1}$$

□

Since this method requires the computation of W_\times , the kernel value can be computed for any type of labeling. For unlabeled graphs, the complexity is $O(kn^3)$ and $O(kdn^3)$ for labeled graphs, where d is the number of labels and k the number of iterations which can be estimated by

$$k = O\left(\frac{\ln \epsilon}{\ln \lambda + \ln |\xi|}\right) \quad (25)$$

Finally, the fixed point iterations is more demanding than the others since it has a condition on the value of λ and thus assumes as geometric $\mu(k)$ function, and the number of iterations can easily vary a lot. The only advantage of this method is that it is very simple while still being better than the raw kernel.

Spectral Decomposition The Spectral Decomposition, usually called the eigendecomposition of a matrix M is a factorization in the form VDV^{-1} where D is the diagonal matrix of eigenvalues and V the matrix of eigenvectors. This decomposition can reduce the computation time of the kernel by making most of the matrix operations trivial

$$\kappa(G, G_2) = \sum_{k=0}^{\infty} \mu(k) q_\times^\top (V_\times D_\times V_\times^{-1})^k p_\times = q_\times^\top V_\times \left(\sum_{k=0}^{\infty} \mu(k) D_\times^k \right) V_\times^{-1} p_\times \quad (26)$$

Indeed, the power now only applies to a diagonal matrix which can be computed much faster. Moreover, this method unlike the others preserves the prior $\mu(k)$ without any conditions. If $\mu(k) = \lambda^k$ is chosen then the kernel becomes

$$\kappa(G, G_2) = q_\times^\top V_\times (I - \lambda D_\times)^{-1} V_\times^{-1} p_\times \quad (27)$$

That kernel is even more trivial to compute as the inverse of a diagonal matrix is simply the inverse of each of its entries and is thus computed in linear time. Finally, by using $\mu(k) = \frac{\lambda^k}{k!}$ as prior, the following new kernel is obtained

$$\kappa(G, G_2) = q_{\times}^{\top} V_{\times} e^{\lambda D_{\times}} V_{\times}^{-1} p_{\times} \quad (28)$$

However, the complexity of the computation of the eigendecomposition is as large as the raw kernel $O(n^6)$, but, for unlabeled graphs, a big improvement is possible by simply calculating the eigendecompositions of the original adjacency matrices A_1 and A_2 thus cutting the time to $O(n^3)$, thanks to a property of the Kronecker product

$$A_1 \otimes A_2 = (V_1 D_1 V_1^{-1}) \otimes (V_2 D_2 V_2^{-1}) = (V_1 \otimes V_2)(D_1 \otimes D_2)(V_1 \otimes V_2)^{-1} \quad (29)$$

By rewriting the spectral decomposition kernel taking advantage of this property, the following new kernel is obtained

$$\kappa(G_1, G_2) = (q_1^{\top} V_1 \otimes q_2^{\top} V_2) \left(\sum_{k=0}^{\infty} \mu(k) (D_1 \otimes D_2)^k \right) (V_1^{-1} p_1^{\top} \otimes V_2^{-1} p_2^{\top}) \quad (30)$$

Which can also be altered by choosing specific $\mu(k)$ as described above. This method ends up being the most efficient of all for unlabeled graphs as it is computed in $O(pn^3)$ where p is the computation time of the power series μ . Moreover this complexity can further improve the computation time of the gram matrix because the eigendecomposition of an adjacency matrix only needs to be computed once. However, this method cannot be expanded to labeled graphs as it encounters difficulties applying the property mentioned above.

Nearest Kronecker Product Approximation This method is more particular as it is not a new way to compute the kernel, but rather a way to generalize all the other methods. Indeed the Nearest Kronecker Product Approximation[18] is used to approximate two matrices A and B as if their Kronecker product was equal to the Kronecker product of two labeled graphs' adjacency matrices $W_{\times} \approx A \otimes B$. This methods allows using any of the methods mentioned above since it returns only two matrices and thus brings back the problem to unlabeled graphs. This method tries to minimize the Frobenius norm of the difference between the two $\|W_{\times} - A \otimes B\|_F$. It can be computed in $O(dn^2)$ since W_{\times} is a sum of Kronecker products.

However this method also brings its own set of drawbacks. Since it required the Kronecker product of each couple of adjacency matrices to be computed, the previous advantage of the Spectral Decomposition (only computing once the eigendecomposition for each adjacency matrix) is lost as the obtained matrices are unique for each W_{\times} .

2.2.4 Digest

tableau cocher avantages inconvenients

remplir

3 Experiments

A big part of this project consisted in implementing the different methods used to compute the random walk kernel[14], verifying the results obtained in that article, and making new attempts to accelerate those algorithms while keeping the best accuracy possible.

3.1 Databases and Metrics

3.1.1 A synthetic graph database

In order to conduct experiments and verify claims made in the article, toy data was necessary. The first thing that had to be done was to create a graph database generator, of very simple and standard graphs that had expectable behaviors and could thus be used for standard tests. Those graphs would also be altered using simple transformations, in order to create different graphs belonging to the same class. For simplicity, it was decided to only use 3 classes to stay as

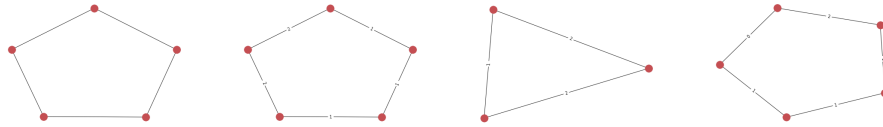


Figure 4: Ring graphs, resp. unlabelled, labelled, altered structured, altered labels

simple as possible. Focus was put on 3 standard types of graphs : the ring, the star and the tree. A ring graph is a connected graph where each vertex is connected to exactly two vertices, it thus forms one big cycle. A star graph is composed of a central node, and of several other nodes that are all and only connected to the central node, it should look like a star, or perhaps a flower.. Finally, a tree is a famous type of graph : it is a connected graph without any cycles, but here a special type of tree will be studying, since each vertex is at most of degree 3.

Est-ce que je devrais dire ça ?

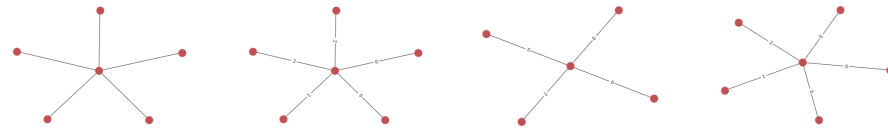


Figure 5: Star graphs, resp. unlabelled, labelled, altered structured, altered labels

graph, it will create a graph of predetermined size with randomly generated edge labels in a certain given interval. Then, for each graph, it will alter it a predefined number of times, and each time will create a new graph that has been both altered on structure and on labels. Alteration on structure involves removing or adding a predefined number of nodes in respect of the type of the graph thus staying in the same class as the source.

1. Ring graphs are simply regenerated slightly longer or shorter randomly.
2. Star graphs are also regenerated either with more or less nodes, randomly.
3. Tree graphs are either expanded by adding randomly leaves anywhere in the graph, or reduced by randomly removing leaves.

Alteration on labels will randomly switch a predefined number of edge labels. Once this is done, the generator will create two databases out of the set of previously generated graphs, one will be made of the adjacency matrix of each graph without looking at labels, and the second one will be an array of adjacency matrices taken from graphs induced by selecting only edges with a specific label, and that for all labels in the label set.

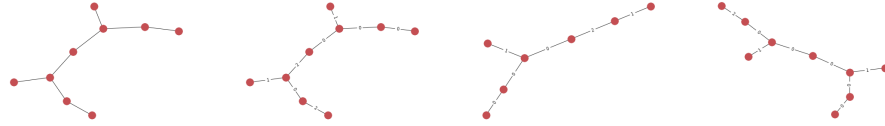


Figure 6: Tree graphs, resp. unlabeled, labeled, altered structured, altered labels

This generator will create databases with the following statistics

# of vertices asked	X	Y
# of vertices per graph		
# of edges per graph		
mean vertex degree		

parler
quand stats
ajoutées

Table 1: Statistics on the synthetic database

3.1.2 Metrics

During the following experiments only two simple metrics will be used to compare the different methods. First, the structural risk referred to as accuracy is

computing simply using the 0-1 loss function as follows

$$L(X, \mathbf{y}) = \frac{1}{|X|} \sum_{i=1}^{|X|} \begin{cases} 1 & \text{if } f(\mathbf{x}_i) \neq y_i \\ 0 & \text{otherwise} \end{cases} \quad (31)$$

This function gives the rate of prediction errors on test data.

The second metric used is the computation time, expressed in seconds.

3.2 Evaluation Protocols

Implementing code from the main paper was a challenge in itself. Several problems aren't really documented and finding functions to solve them, or even explanations on how to solve them is sometimes difficult.

Support Vector Machines The Support Vector Machines algorithm was not implemented since it was not the focus of this project. Instead, the famous library `libsvm`[19] was used through the python library `scikit-learn`[20] because of its efficiency and renown. Moreover, there are some hyperparameters to optimize, namely the penalty coefficient on classification mistakes C . Since it is not the topic of this project it has not been optimized.

Faire ça
rapidement

Gram Matrix todo

faire ça

Raw Kernel The Raw Kernel is very simple to implement however encounters an issue : a sum to infinity cannot be computed as is, and thus a choice had to be made. It was decided that the raw kernel would be computed over N iterations where N would be the number of vertices of the largest graph. The advantage is that all graph kernel values would be computed with the same number of iterations and the number of iterations would be kept small while all vertices would be reachable from any other vertex (the shortest path between two vertices in a connected graph $G = (V, E)$ is at most $|V|$ long).

Inverse Kernel This kernel was also very simple to implement since it only required computing $I - \lambda W_\times$ and its inverse, which was handled exclusively by numpy (part of the scipy project).

Sylvester Equation The Sylvester Equation is a good example of poorly documented problems (compared to the others). There are very few resources on the subject, and even less implementation of solvers. We could theoretically use the Sylvester Equation on both labeled and discrete-labeled graphs which are respectively represented as Sylvester and Generalized Sylvester Equations (involving a sum of matrices), however, the only library available solves the first one, but not the second one. There seems to be confusing terminology since there is a solver for the Generalized Sylvester Equation, but a different one than the one described in the studied source[14]. The same source also

mentions a way to solve the generalized version[21], however, this option was not explored because of a lack of time.

changé ça ?

Conjugate Gradient Method The Conjugate Gradient Method is however very well documented[17] and there are several libraries implementing this algorithm. In this case, we first used the function *cg* of scipy[22] but the results weren't satisfying. After some research, it was found that providing a preconditioner for the inverse of the matrix $I - \lambda W_{\times}$ in the form of a random vector \mathbf{x} multiplied by itself $\mathbf{x} \cdot \mathbf{x}^{\top}$ largely improved the results. Moreover the algorithm requires the matrix to be positive definite thus the matrix was symmetrized. However this function only handles sparse matrices and gave very poor results. It was then decided to implement it directly and gave much better results. As said previously, the algorithm guarantees convergence in n step where n is the effective rank of the matrix $I - \lambda W_{\times}$. However, by grid search on the max-



Figure 7: Accuracy and computation time depending on number of maximum iterations

imum number of iterations allowed for the algorithm, it was discovered that on synthetic data the effective rank of the matrices are very low and the algorithm converges very quickly. This can be cause by the fact that the synthetic database is too simplistic and the adjacency matrices too sparse, making the problem too easy to solve.

Fixed Point Iterations The Fixed Point Iterations method was very easy to implement since the algorithm only consists in repeating $x_{t+1} = p_{\times} + \lambda W_{\times} x_t$ until either convergence is obtained or the maximum number of iterations reached. It rapidly gave good results thus being used as a tool of comparison with other kernels that struggled more at the beginning. However there is a constraint on the lambda (the geometric factor) used in the kernel as proved above. Experiments showed that a lambda very close to the constraint increases significantly

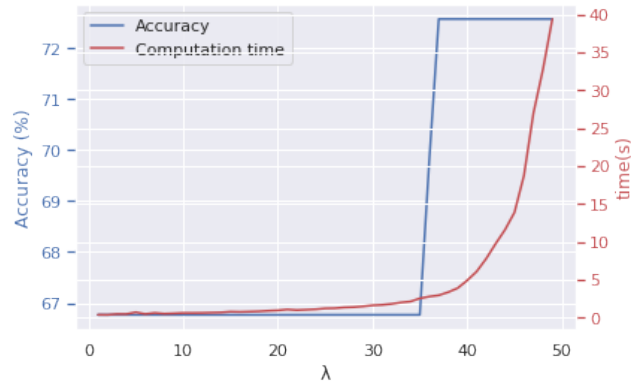


Figure 8: Accuracy and computation time depending on chosen λ

the time of computation to convergence, as well as the density of the adjacency matrix to a lower extent.

Spectral Decomposition The Spectral Decomposition method is also very well documented since it is simply another name for the eigendecomposition, one of the main problems of linear algebra. Because of that it was very simple to implement while following all the smaller accelerations recommended. However, in some cases on the synthetic database the eigenvalues would not be real thus causing some problems. In order to prevent that, the adjacency matrices were also symmetrized. Moreover, in some other cases the left eigenvectors matrix was singular and was not invertible. In order to fix that, an ϵ was added to all the diagonal such that it wouldn't alter the results (and didn't) while allowing all kernels to be computed.

Nearest Kronecker Product En travail

finir

[https://www.sciencedirect.com/science/article/pii/S0377042700003939?](https://www.sciencedirect.com/science/article/pii/S0377042700003939?via=ihub)

[https://www.imageclef.org/system/files/CLEF2016_Kronecker_Decomposition.](https://www.imageclef.org/system/files/CLEF2016_Kronecker_Decomposition.pdf)

http://dx.doi.org/10.1007/978-94-015-8196-7_17

<https://mathematica.stackexchange.com/questions/91651/nearest-kronecker-product>

3.3 Performance

First of all the different methods will be tested with synthetic data to verify claims made in the article. Then, the methods will be applied to real world problems.

3.3.1 Performance on synthetic data

Gram Matrices The algorithms approximate the kernel using very different techniques, thus the gram matrices obtained from the kernels were similar in appearance, but had significant differences of scale. It was then decided to normalize very simply the gram matrices so they could be more easily compared.

$$M = (M - \min(M)) / \max(M) \quad (32)$$

It can be seen that the two matrices are indeed very similar, but locally there



Figure 9: A gram matrix computed with the raw method, another with the fixed point method, and the difference between the two

can be big differences. The difference matrix shows that the fixed point kernel has higher values on the three diagonal blocks in average, and in average also, less in the outer parts of the matrix. Since in this example the data is ordered by class, it shows that the fixed point method gives in average higher values to graphs of the same class and smaller values to graphs of different classes than the raw method. It can be conjecture that the fixed point method will give a better accuracy, which is verified empirically in another paragraph. Afterwards, in order to verify all the algorithms gave similar results, it was decided to compute the Frobenius norm of the difference of two gram matrices (divided by the size of the matrix, thus giving the mean standard deviation). The following results were obtained and were satisfying. Indeed, the mean standard deviation

	Raw kernel	Inverse Kernel	Sylvester Equation	Conjugate Gradients	Fixed points	Spectral Decomp.
Raw.	0	1.1e-4	9.8e-5	8.9e-5	1.0e-4	1.0e-04
Inv.	-	0	2.1e-5	7.9e-5	4.0e-6	6.8e-6
Syl.	-	-	0	8.0e-5	1.7e-5	1.4e-5
Con.	-	-	-	0	7.9e-5	7.9e-5
Fix.	-	-	-	-	0	2.8e-6
Spe.	-	-	-	-	-	0

Table 2: Mean standard deviation of matrix entries

is significantly enough small in most cases so that the different acceleration methods are probably going to give approximately correct results.

Label Use The synthetic database generates at the same time adjacency matrices of each graph, and the list of adjacency matrices filtered by each label. Thus, we can compare the scores obtained for each of our metrics, computed using the fixed point kernel since it is both fast and allows the use of labels. The accuracy has a lot of variance and thus it is hard to compare it. It can

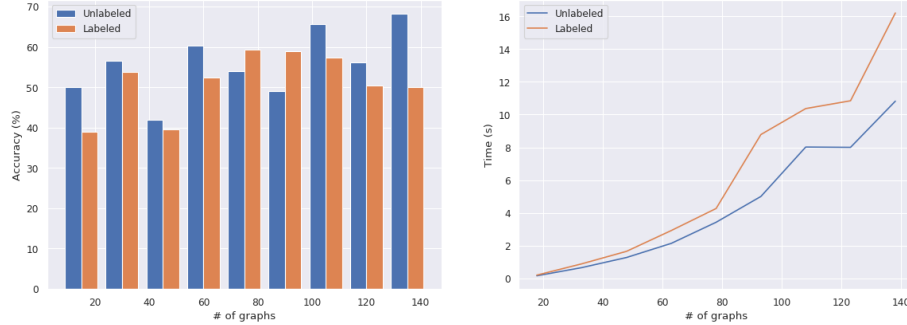


Figure 10: Accuracy and computation time of learning for unlabeled and labeled graphs

probably be explained by the fact that we are working on graphs that are not big enough, and thus a small change in labels could have a big effect on the kernel, for example on the central edge of a tree graph. It is conjectured that provided a graph large and connected enough that a small change in label does not alter much its list of adjacency matrices, labels will have a positive impact on the accuracy. Concerning the computation time, the difference is quite small though by increasing the size of the database the difference can be better observed. There is a small factor of difference during the computation of the tensor product.

refaire cette
experi-
ence en
changeant
quoi ?

Comparison of methods depending on number of graphs Pas de sens, grandit à la meme vitesse = taille de bd au carré ?

Comparison of methods depending on number of nodes In order to verify the complexity of the different methods and compare their efficiency their accuracy and computation time was calculated on synthetic databases with graphs of varying sizes. The databases were generated independently so the different entries cannot be compared exactly, but being generated following the same precise rules, the deviation is not too large. Concerning the accuracy, the conjecture made earlier is verified in most cases : the fixed point method, except on the largest graphs, gives among the best scores. It can be also noted that the raw kernel method is rarely beaten, only twice by the fixed point kernel only. Surprising the inverse method performs in average worst, all the other methods approximating it giving better results.

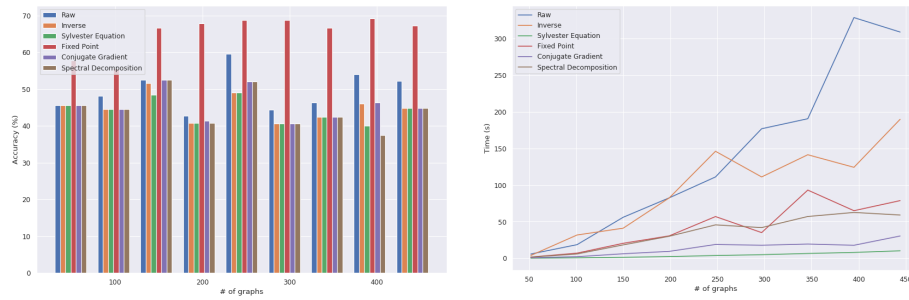


Figure 11: Accuracy and Computation time of different methods depending on the size of the database

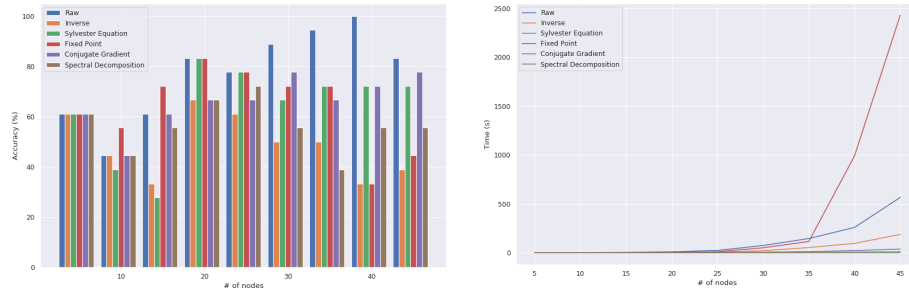


Figure 12: Accuracy and Computation time of different methods depending on the size of graphs

On the other hand, computation time. It is important to remind that in this case some methods (fixed point, conjugate gradient) were implemented entirely for this project, while the Spectral Decomposition and Inverse methods used tools to compute eigenvalues and inverses, and finally the Sylvester Equation solver was taken from a library based in Fortran. Fortran being much faster than python, it creates a bias in this graph, however it is not an issue since we will only look at the evolution of curves, not their actual values.

Nearest Kronecker product En implémentation

3.3.2 Performance on real datasets

This implementation was also tested on two real databases being PROTEINS and ENZYMES TABLEAU

arranger
point fixe
et finir

finir

4 Conclusion and Future Work

experiences : détailler db, tests, méthodes, les parametres, construction label,
donner tableaux, resultats, teps de calculs, précision, discuter tout cela section
4 : 1 page ou page et demi : conclusion et discussion

A Appendix

Proofs ?

Acknowledgments

This work was done during the first year of my master's at Sorbonne Université. I would really like to thank Professor Sahbi for all the time he dedicated to me, and for his extraordinary and expert coaching on this project. This project taught me a lot and made me closer to becoming a researcher, while also making me a more mature person.

Mettre ces
listes ou pas
?

List of Figures

1	Two "tree" graphs, resp. unlabeled and labeled	4
2	A graph G_1 , a graph G_2 and the product graph $G_1 \otimes G_2$	9
3	Conjugate gradient (red) compared to Gradient Descent (green)	11
4	Ring graphs, resp. unlabelled, labelled, altered structured, altered labels	15
5	Star graphs, resp. unlabelled, labelled, altered structured, altered labels	15
6	Tree graphs, resp. unlabeled, labeled, altered structured, altered labels	16
7	Accuracy and computation time depending on number of maximum iterations	18
8	Accuracy and computation time depending on chosen λ	19
9	A gram matrix computed with the raw method, another with the fixed point method, and the difference between the two	20
10	Accuracy and computation time of learning for unlabeled and labeled graphs	21
11	Accuracy and Computation time of different methods depending on the size of the database	22
12	Accuracy and Computation time of different methods depending on the size of graphs	22

List of Tables

1	Statistics on the synthetic database	16
2	Mean standard deviation of matrix entries	20
bibliographie et index		

References

- [1] J. A. Bondy, U. S. R. Murty, *et al.*, *Graph theory with applications*, vol. 290. Citeseer, 1976.
- [2] K. M. Borgwardt, C. S. Ong, S. Schönauer, S. V. N. Vishwanathan, A. J. Smola, and H.-P. Kriegel, “Protein function prediction via graph kernels,” *Bioinformatics (Oxford, England)*, vol. 21 Suppl 1, pp. i47–56, June 2005.
- [3] C. J. C. Burges, “A Tutorial on Support Vector Machines for Pattern Recognition,” *Data Min. Knowl. Discov.*, vol. 2, pp. 121–167, June 1998.
- [4] C. Cortes and V. Vapnik, “Support-vector networks,” *Machine Learning*, vol. 20, pp. 273–297, Sept. 1995.
- [5] Y. Freund and R. E. Schapire, “Large margin classification using the perceptron algorithm,” *Machine learning*, vol. 37, no. 3, pp. 277–296, 1999.
- [6] D. Haussler, “Convolution kernels on discrete structures,” Technical Report UCS-CRL-99-10, University of California at Santa Cruz, Santa Cruz, CA, USA, 1999.
- [7] H. Kashima, K. Tsuda, and A. Inokuchi, “Marginalized kernels between labeled graphs,” in *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, ICML’03, pp. 321–328, AAAI Press, 2003.
- [8] N. Shervashidze, S. V. N. Vishwanathan, T. Petri, K. Mehlhorn, and K. Borgwardt, “Efficient graphlet kernels for large graph comparison,” in *Artificial Intelligence and Statistics*, pp. 488–495, Apr. 2009.
- [9] J. Ramon and T. Gärtner, “Expressivity versus efficiency of graph kernels,” Citeseer.
- [10] K. M. Borgwardt and H.-P. Kriegel, “Shortest-path kernels on graphs,” in *Fifth IEEE international conference on data mining (ICDM’05)*, pp. 8–pp, IEEE, 2005.
- [11] N. Shervashidze, *Scalable graph kernels*. PhD thesis, Universität Tübingen, 2012.
- [12] B. Schölkopf and A. J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, Mass: The MIT Press, 1st edition ed., Dec. 2001.
- [13] T. Weissman, E. Ordentlich, G. Seroussi, S. Verdu, and M. J. Weinberger, “Inequalities for the l_1 deviation of the empirical distribution,” *Hewlett-Packard Labs, Tech. Rep.*, 2003.

- [14] S. V. N. Vishwanathan, N. N. Schraudolph, R. Kondor, and K. M. Borgwardt, “Graph Kernels,” *Journal of Machine Learning Research*, vol. 11, no. Apr, pp. 1201–1242, 2010.
- [15] W. Imrich and S. Klavzar, *Product graphs: structure and recognition*. Wiley, 2000.
- [16] P. M. Weichsel, “The kronecker product of graphs,” *Proceedings of the American mathematical society*, vol. 13, no. 1, pp. 47–52, 1962.
- [17] Y. Nesterov, *Lectures on Convex Optimization*. Springer Optimization and Its Applications, Springer International Publishing, 2 ed., 2018.
- [18] C. F. Van Loan and N. Pitsianis, “Approximation with kronecker products,” in *Linear algebra for large scale and real-time applications*, pp. 293–314, Springer, 1993.
- [19] C.-C. Chang and C.-J. Lin, “LIBSVM: A library for support vector machines,” *ACM Transactions on Intelligent Systems and Technology*, vol. 2, pp. 27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [20] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, *et al.*, “Scikit-learn: Machine learning in python,” *Journal of machine learning research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [21] L. De Lathauwer, B. De Moor, and J. Vandewalle, “Computation of the canonical decomposition by means of a simultaneous generalized schur decomposition,” *SIAM Journal on Matrix Analysis and Applications*, vol. 26, no. 2, pp. 295–327, 2004.
- [22] E. Jones, T. Oliphant, P. Peterson, *et al.*, “Scipy: Open source scientific tools for python, 2001,” 2016.