

Learning to Aggregate on Structured Data

Clemens Damke

November 9, 2019



PADERBORN UNIVERSITY
The University for the Information Society

Department of Electrical Engineering,
Computer Science and Mathematics
Warburger Straße 100
33098 Paderborn



INTELLIGENT
SYSTEMS

Intelligent Systems Group (ISG)

Master Thesis

Learning to Aggregate on Structured Data

Clemens Damke

- | | |
|-------------|--|
| 1. Reviewer | Prof. Dr. Eyke Hüllermeier
Intelligent Systems and Machine Learning Group (ISG)
Paderborn University |
| 2. Reviewer | Prof. Dr. Axel-Cyrille Ngonga Ngomo
Data Science Group (DICE)
Paderborn University |

November 9, 2019

Clemens Damke

Learning to Aggregate on Structured Data

Master Thesis, November 9, 2019

Reviewers: Prof. Dr. Eyke Hüllermeier and Prof. Dr. Axel-Cyrille Ngonga Ngomo

Supervisor: Vitalik Melnikov

Paderborn University

Intelligent Systems and Machine Learning Group (ISG)

Heinz Nixdorf Institute

Department of Electrical Engineering, Computer Science and Mathematics

Warburger Straße 100

33098 Paderborn

Abstract

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Goals	2
1.3	Structure	2
2	Related Work	3
2.1	Learning to Aggregate	3
2.2	Graph Regression and Classification	4
2.2.1	Graph Embeddings	4
2.2.2	Graph Kernels	4
2.2.3	Graph Neural Networks	4
3	Learning to Aggregate on Graphs	5
3.1	Formalization of LTA Methods	5
3.2	An LTA Interpretation of Graph Methods	5
3.3	LTA on Dynamically Decomposed Graphs	5
4	Evaluation	7
5	Conclusion	9
5.1	Review	9
5.2	Future Directions	9
A	Appendix	13
	Bibliography	15

Introduction

1.1 Motivation

The field of *machine learning* (ML) on graph-structured data has applications in many domains due to the general expressive power of graphs. Three common types of graph ML problems are

1. **Link prediction:** A graph with an incomplete edge set is given and the missing edges have to be predicted. The generation of friendship suggestions in a social network is a typical example for this.
2. **Vertex classification & regression:** Here a class or a score has to be predicted for each vertex of a graph. In social graphs this corresponds to the prediction of properties of individuals. Another example is the prediction of the amount of traffic at the intersections of a street network.
3. **Graph classification & regression:** In this final problem type a single global class or continuous value has to be predicted for an input graph. The canonical example for this is the prediction of properties of molecule graphs, e.g. the toxicity or solubility of a chemical.

In this thesis we will focus on the last problem type, *graph classification and regression* (GC/GR). A ML method for this problem has to accept graphs of varying size and should be permutation invariant wrt. the vertices. Those requirements are not met by the commonly used learners that only accept fixed-size feature vectors as their input, e.g. logistic models, support vector machines or multilayer perceptrons.

A GC/GR method has to account for two central aspects of the problem: 1. Local structural analysis and 2. global aggregation. The first aspect is about the extraction of relevant features of substructures of the input graph. The latter is about the way in which the local features are combined into a final class or regression value. The existing GC/GR methods are mostly motivated by local structural graph analysis. The aspect of global aggregation on the other hand is less emphasized by those methods.

There is however a separate branch of research that specifically looks at the problem of learning aggregation functions, called *learning to aggregate* (LTA). Current LTA approaches explicitly learn an aggregation functions for sets which can be interpreted

as graphs without edges. The motivation for this thesis is to generalize LTA from sets to arbitrary graphs. The overall goal is to combine the aggregation learning perspective with existing GC/GR methods.

1.2 Goals

To extend LTA to graphs, three goals have to be achieved:

1. **Formalization of LTA:** Before LTA can be extended, its essential characteristics have to be defined. Those characteristics should provide the terminology to formally capture the differences and similarities between LTA and existing GC/GR methods.
2. **Give an LTA interpretation of GC/GR methods:** Using the LTA formalization, representative GC/GR approaches should be restated as LTA instances. Currently there is no comprehensive formulation of the relation between both fields of research; this is addressed by the the second goal.
3. **Define an LTA method for graphs:** Using the LTA perspective on GC/GR, hidden assumptions of the existing approaches should become clear and in which way they share the assumptions of LTA. The last goal is to use those insights to formulate an LTA-GC/GR method that combines ideas from the existing approaches with the LTA assumptions.

1.3 Structure

Chapter 2: Related Work

Chapter 3: Learning to Aggregate on Graphs

Chapter 4: Evaluation

Chapter 5: Conclusion

Related Work

Before combining LTA and GC/GR as described in section 1.2, we first give an overview of the state-of-the-art in both fields of research. First the existing LTA methods for unstructured set inputs are described. Then we will look at the two main families of GC/GR methods: 1. Vector representation methods 2. Graph neural networks.

2.1 Learning to Aggregate

The class of LTA problems was first described by Melnikov and Hüllermeier [MH16]. There an input instance is understood as a composition c_i of so-called constituents $c_{i,j} \in c_i$, i.e. as a variable-size multiset with $n_i = |c_i|$. The assumption in LTA problems is that for all constituents $c_{i,j}$ a local score $y_{i,j} \in \mathcal{Y}$ is either given or computable. The set of those local scores should be indicative of the overall score $y_i \in \mathcal{Y}$ of the composition c_i . LTA problems typically require two subproblems to be solved:

1. **Aggregation:** A variadic aggregation function $A : \mathcal{Y}^* \rightarrow \mathcal{Y}$ that estimates composite scores has to be learned, i.e. $y_i \approx \hat{y}_i = A(y_{i,1}, \dots, y_{i,n_i})$. Typically the aggregation function A should be associative and commutative to fit with the multiset-structure of compositions.
2. **Disaggregation:** In case the constituent scores $y_{i,j}$ are not given, they have to be derived from a constituent representation, e.g. a vector $x_{i,j} \in \mathcal{X}$. To learn this derivation function $f : \mathcal{X} \rightarrow \mathcal{Y}$, only the constituent vectors $\{x_{i,j}\}_{j=1}^{n_i}$ and the composite score y_i is given. Thus the constituent scores $y_{i,j}$ need to be *disaggregated* from y_i in order to learn f .

Overall LTA can be understood as the joint problem of learning the aggregation function A and the local score derivation function f . Two main approaches to represent the aggregation function in LTA problems have been explored.

Uninorm-Aggregation The first approach uses *uninorms* [MH16] to do so. There the basic idea is to express composite scores as fuzzy truth assignments $y_i \in [0, 1]$. Such a composite assignment y_i is modeled as the result of a parameterized logical expression of constituent assignments $y_{i,j} \in [0, 1]$. As the logical expression that thus effectively aggregates the constituents, a uninorm U_λ is used. Depending on the parameter λ , U_λ combines t-norms and t-conorms which are continuous generalizations of logical conjunction and disjunction respectively.

OWA-Aggregation Recently Melnikov and Hüllermeier [MH19] have also looked at an alternative class of aggregation functions. Instead of using fuzzy logic to describe score aggregation, *ordered weighted average* (OWA) operators were used. OWA aggregators work by sorting the input scores and then weighting them based on their sort position, i.e.

$$A_{\lambda}(y_1, \dots, y_n) := \sum_{i=1}^n \lambda_i y_{\pi(i)},$$

where λ is a weight vector with $\|\lambda\|_1 = 1$ and π is a sorting permutation of the input scores. To deal with varying composition sizes n the weights λ_i are interpolated using a *basic unit interval monotone* (BUM) function $q : [0, 1] \rightarrow [0, 1]$. It takes constituent positions that are normalized to the unit interval, i.e. $\frac{i}{n}$. The BUM function q then is then used to interpolate a weight for any normalized sort position via $\lambda_i = q(\frac{i}{n}) - q(\frac{i-1}{n})$. Therefore the learning objective boils down to optimizing the shape of q .

2.2 Graph Regression and Classification

2.2.1 Graph Embeddings

2.2.2 Graph Kernels

2.2.3 Graph Neural Networks

Spatial GNNs

Spectral GNNs

Learning to Aggregate on Graphs

3.1 Formalization of LTA Methods

3.2 An LTA Interpretation of Graph Methods

3.3 LTA on Dynamically Decomposed Graphs

Conclusion

5.1 Review

5.2 Future Directions

Appendix

A

Bibliography

- [MH16] Vitalik Melnikov and Eyke Hüllermeier. „Learning to Aggregate Using Uninorms“. In: Machine Learning and Knowledge Discovery in Databases. Springer International Publishing, 2016, pp. 756–771 (cit. on p. 3).
- [MH19] Vitalik Melnikov and Eyke Hüllermeier. „Learning to Aggregate: Tackling the Aggregation/Disaggregation Problem for OWA“. In: ACML (2019) (cit. on p. 4).

List of Figures

Erklärung zur Masterarbeit

Ich, Clemens Damke (Matrikel-Nr. 7011488), versichere, dass ich die Masterarbeit mit dem Thema *Learning to Aggregate on Structured Data* selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Die Stellen der Arbeit, die ich anderen Werken dem Wortlaut oder dem Sinn nach entnommen habe, wurden in jedem Fall unter Angabe der Quellen der Entlehnung kenntlich gemacht. Das Gleiche gilt auch für Tabellen, Skizzen, Zeichnungen, bildliche Darstellungen usw. Die Masterarbeit habe ich nicht, auch nicht auszugsweise, für eine andere abgeschlossene Prüfung angefertigt. Auf § 63 Abs. 5 HZG wird hingewiesen.

Paderborn, 9. November 2019

Clemens Damke