

Jelena Tešić: Research Impact and Plan

Unstructured data encompasses diverse data types and formats, such as video data, electronic health data records, surveys, experimental readings, sensor data, social platform content data, genomics data, and purchase reviews. The arching quest of my research is to provide practical, efficient, effective, intuitive, and responsible algorithms for addressing the challenges of analytic tasks when applied to such extensive unstructured data collections in the wild. The research I lead in the Data Lab @ TXST (DataLab12.github.io) is grouped in the following sections by the tasks it solves.

1 Scaling Signed Graph Balancing Tasks through Fundamental Cycle Basis

In network science, signed graph representation of relations provides more information than unsigned graph networks. Spectral pollution, computational complexity, inherent bias, and training requirements are the main hurdles that have slowed the adoption of the signed network as data representation and signed graph algorithms for modern data analysis. We took a different approach, stemming from mathematical sociology, and focused on balance theory to solve the signed network tasks at scale. First, we have expanded the balance theory to signed social network graph analysis. We propose a frustration cloud view of the signed graph where we quantified vertex and edge in terms of frustration cloud statistics and validated this novel social network graphs analysis approach [30] [pdf]. Next, we have developed an algorithm to efficiently compute the fundamental cycle basis in large, unstructured graphs without requiring training data or relying on spectral computation assumptions to scale the findings to signed graphs constructed from social networks and recommendation data [1]. There was no benchmark to evaluate state-of-the-art tasks on signed graphs derived from real networks, so we have created the most extensive assumption-free comparison of community discovery on signed networks in [14] [pdf] in terms of efficacy, efficiency, scalability, and reproducibility of existing methods, and propose frustration cloud-based approach for cluster boosting for high modular signed graphs [15] [pdf]. Two master and five undergraduate students contributed to this research project.

The research plan is to use a fundamental cycle basis to scale the solution of other NP-hard tasks to large signed networks using the efficient fundamental cycle computation approach [22] [pdf] and propose algorithms to solve the task of computing frustration and the balanced state of the graph at scale [24] [pdf] and the task of finding the largest balanced subgraph in any network [23] [pdf] for graphs with millions of nodes and edges. The project has evolved into a significant part of the Ph.D. thesis work. Next, we plan to propose a scalable community discovery algorithm and extend the analysis using a fundamental cycle basis to the recommendation and anomaly detection tasks, as well as to extend the application to sensor and agent networks and signed gene networks.

2 Modeling Social Network Relations and Improving Label Propagation

The multifaceted interconnectivity of users and content on Twitter through user connections, replies, quotes, hashtags, and shared content makes it an exciting medium for research on the effectiveness of the representation and methods used. First, we have introduced a scalable end-to-end Twitter network data management pipeline that gathers, stores, and models rich relationships from Twitter networks [12] [pdf]. Next, we compared and contrasted the analysis results of millions of Twitter data using multiple graph construction processing approaches [13] [pdf]. How well the tweet content can be classified based on modeling relationships from interactions alone, and how well can community classification predict the label of the content? The community-based modeling (tweet is classified not on the content but on the retweets, replies, quotes, hashtags, and the author) yields precision, recall, and accuracy comparable to

lexical classifiers [21] [pdf]. We propose novel multi-modal approaches that consistently deliver the most robust outcomes and exhibit the highest performance measures for network graphs constructed based on Twitter interactions related to the COVID-19 pandemic [28] [pdf] and [21] [pdf]. One undergraduate, one master, and one Ph.D. student participated in the research work for this research project. The next step is to extend relationship modeling to extended signed graphs and hyper-graphs; and to explore network representations to provide an assumption-free and bias-free baseline for evaluating graphical neural network performance for network science tasks.

3 Identifying Small Objects in Highly Variable Overhead Videos

The high variability of content in the overhead imagery stems from the terrestrial region captured, the high variability of acquisition conditions, and the number and size of objects in aerial imagery that are very different than in the consumer data. State-of-art fails due to the high variability of the domain and the low availability of training data. We have proposed multiple domain adaptation approaches to alleviate the degradation of object identification in previously unseen overhead datasets with significant domain gaps and dominant small objects [5] [pdf] and [25] [pdf]. We have proposed novel algorithms for overhead videos to detect anomalous activities [9, 32]. Five undergraduates, two masters, and three Ph.D. students participated in this project sponsored by NAVAIR. The project evolved in the Ph.D. thesis work, where we have proposed multiple innovative and efficient contrastive learning algorithms to improve object classification in previously unseen highly variable overhead datasets [6] [pdf] and [8] [pdf]. We are now introducing progressive domain adaptation to produce domain-invariant features across aerial datasets using local and global components for domain adaptation and object classification for the task [7] [pdf] and [17] [pdf]. We will focus on improving scene understanding in highly variable videos.

4 Deep Descriptor Database Indexing Search and Retrieval

Searching for unseen objects in extensive visual archives is challenging, demanding efficient indexing methods supporting meaningful similarity retrievals. Together with the Ph.D. student, I have introduced novel indexing and search algorithm for deep descriptor databases that have up to four times lower memory usage and higher effectiveness than state-of-art [19] [pdf] and [20] [pdf] on millions of deep descriptors. We are building upon the approach for crowd-sensing application [17] [pdf] and plan to improve the indexing footprint while keeping the search effectiveness by proposing a novel stratified graph approach [18] [pdf]. Next, we will focus on multi-modal indexing, combining feature indexing with GIS attribute, spatial, and shape indexing.

5 Semantic Segmentation Task in The Wild

Semantic image segmentation is the task of assigning a label to each pixel in an image. The task has various practical applications in medicine, transportation, machine vision, and science. My contribution to semantic segmentation focuses on pavement distress detection for transportation and road maintenance applications. We have quantified the main influencing factors that affect the performance of deep learning models in pavement distress detection pipelines and proposed a semantic segmentation algorithm that significantly improves the accuracy of localizing pavement cracks [27] [pdf]. Next, we explore the domain adaptation approaches from section 3. to improve the algorithmic performance for this specialized field. We are also designing a new semantic segmentation pipeline to automatically classify Bacterial Adhesion and Corrosion from images obtained in the NASA SpaceX-21 experiment. The inability to control microbial biofilms during spaceflight poses a severe health risk to astronauts. Automatically identifying corrosion in tens of thousands of images of samples flown into space will help researchers streamline the samples'

data collection and help them reach conclusions faster on what countermeasures work. Two undergraduates, two masters, and one Ph.D. student are actively involved in this project.

6 Predictive Modeling of Noisy Tabular Data

Tabular data in the wild are difficult to model due to the uneven distribution of attributes, missing, overlapping, noisy values, a mix of categorical and numerical data attributes, data imbalance, and a long tail of sparse values. We have developed the intentional data science pipeline that can automatically uncover important attributes, reduce feature space, and model prediction in a robust manner from multi-source tabular data in [10]. We have designed the multi-feature importance analysis algorithm and applied it to large-scale analysis of public data from the National Center for Education Statistics (NCES) to provide data-driven insights into teacher attrition challenges. We discovered that the race and sex of the principal, the type of school, and the school's location impact teacher retention rates the most and that modeling historical data resulted in a predicted attrition rate of over 10%, aligning closely with the current prevalent attrition rates in the USA [11] [pdf]. Next, we have developed an interpretable data-driven scoring fusion to discover the most critical factors from an extensive collection of heterogeneous public data sources on learning loss during the COVID-19 pandemic in Texas public schools. Our robust approach found that the number of students on school campuses in the Fall of 2020 and in the Spring of 2021 was the most resilient and most impactful predictor of how the students would perform on the standardized test in mathematics and reading in the Spring of 2021 in Texas [4] [pdf]. Finally, we introduced novel cascade enhancement to ensure effectiveness and the prediction coverage of our modeling pipeline to predict long COVID in N3C data [16] [pdf]. The project resulted in one honors thesis [3], one master thesis [10], and as a seed to a Ph.D. thesis work [16] [pdf]. Next, we will improve the loss function in our gradient boosting approach to handle the sparsity, relevance, and hierarchical values in the tabular N3C data more effectively. We will use temporal modeling and the latest educational data to integrate learning recovery predictions.

7 Summary

Algorithms for unstructured data are driven by the data type, format, acquisition, size, and intended task. My research focuses on providing novel algorithms that derive from the underlying mathematics, computer science, and statistics theory while considering the specifications (efficiency, scalability, usability, interpretability) of the task at hand and the domain applicability. To this end, I have collaborated with other Labs at our department and proposed data-driven solutions for their specific challenges [31, 2, 26, 29]. I am currently looking into developing implicit neural representations to encode large climate models. My plan for Data Lab is to continue to invent novel algorithms and methods for solving challenging analytics tasks stemming from the nature and size of unstructured data.

* *author* - names of the the Data Lab students are in *italic*

References Cited

- [1] Ghadeer Alabandi, **Jelena Tešić**, Lucas Rusnak, and Martin Burtscher. Discovering and balancing fundamental cycles in large signed graphs. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '21, New York, NY, USA, 2021. Association for Computing Machinery.
- [2] Noah Dunstatter, Alireza Tahsini, Mina Guirguis, and **Tešić, Jelena**. Solving cyber alert allocation markov games with deep reinforcement learning. In Tansu Alpcan, Yevgeniy Vorobeychik, John S. Baras, and György Dán, editors, *Decision and Game Theory for Security*, pages 164–183, Cham, 2019. Springer International Publishing.
- [3] *Daniel Payan*. Ranking resilience attributes for texas public school districts. Honor’s thesis, TXST, 2023. Advisor: **Jelena Tešić**.
- [4] *Daniel Payan*, *June Yu*, Li Feng, and **Jelena Tešić**. Data driven intervention for covid learning loss in texas public school districts. *IEEE Transactions on Learning Technologies* (*under review*).
- [5] *David Heyse*, *Nicholas Warren*, and **Jelena Tešić**. Identifying maritime vessels at multiple levels of descriptions using deep features. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, pages 423 – 431. SPIE, 2019.
- [6] *Debojyoti Biswas*, *MMM Rahman*, Ziliang Zong, and **Jelena Tešić**. Improving the energy efficiency of real-time dnn object detection via compression, transfer learning, and scale prediction. In *The IEEE 16th International Conference on Networking, Architecture, and Storage (NAS 2022)*, September 2022.
- [7] *Debojyoti Biswas* and **Jelena Tešić**. Domain adaptation with contrastive learning for object detection in satellite imagery. *IEEE Transactions on Geoscience and Remote Sensing* (*under review*).
- [8] *Debojyoti Biswas* and **Jelena Tešić**. Small object difficulty modeling for objects detection in satellite images. In *IEEE 14th International Conference on Computational Intelligence and Communication Networks (CICN)*, pages 125–130, 2022.
- [9] *George E. Strauch*, *Jiajian (Jax) Lin*, and **Jelena Tešić**. Overhead projection approach for multi-camera vessel activity recognition. In *2021 IEEE International Conference on Big Data (Big Data)*, pages 5626–5632, 2021.
- [10] *June Yu*. Gradient boosting public data modeling for the policy planning in education. Master’s thesis, TXST, 2022. Advisor: **Jelena Tešić**.
- [11] *June Yu*, Li Feng, and **Jelena Tešić**. Mitigating u.s. public school teacher attrition crisis: A data science approach. *Information Processing & Management* (*under review*).
- [12] *Lia Nogueira de Moura*. Social network analysis at scale: Graph-based analysis of twitter trends and communities. Master’s thesis, TXST, 2020. Advisor: **Jelena Tešić**.

- [13] *Lia Nogueira de Moura* and **Jelena Tešić**. pytwanalysis: Twitter data management and analysis at scale. In *2021 Seventh International Conference on Social Networks Analysis, Management and Security (SNAMS)*, 2021.
- [14] *Maria Tomasso*, Lucas Rusnak, and **Jelena Tešić**. Advances in scaling community discovery methods for signed graph networks. *Journal of Complex Networks*, 10(3), 06 2022. cnac013.
- [15] *Maria Tomasso*, Lucas Rusnak, and **Jelena Tešić**. Cluster boosting and data discovery in social networks. In *Proceedings of the 37th ACM/SIGAPP Symposium On Applied Computing (SAC)*, 2022.
- [16] *Mirna Elizondo*, Rasim Musal, *June Yu*, and **Jelena Tešić** on behalf of N3C. Long covid challenge: Predictive modeling of noisy clinical tabular data. In *2023 IEEE 11th International Conference on Healthcare Informatics (ICHI)*, 2023.
- [17] *MMM Rahman*, *Debojyoti Biswas*, and **Jelena Tešić**. Evirec: Efficient visual indexing and retrieval for edge crowd-sensing. In *Submitted to a Conference*.
- [18] *MMM Rahman* and **Jelena Tešić**. Stratified graph indexing for efficient search in deep descriptor databases. In *Submitted to a Conference*.
- [19] *MMM Rahman* and **Jelena Tešić**. Evaluating hybrid approximate nearest neighbor indexing and search (hannis) for high-dimensional image feature search. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 6802–6804, 2022.
- [20] *MMM Rahman* and **Jelena Tešić**. Hybrid approximate nearest neighbor indexing and search (hannis) for large descriptor databases. In *2022 IEEE International Conference on Big Data*, pages 3895–3902, 2022.
- [21] *Muhieddine Shebaro*, *Lia Nogueira de Moura*, and **Jelena Tešić**. Multimodal mining of twitter networks for improved label propagation. *Social Network Analysis and Mining (under review)*.
- [22] *Muhieddine Shebaro* and **Jelena Tešić**. Identifying stable states of large signed graphs. In *Companion Proceedings of the ACM Web Conference 2023 (WWW '23 Companion)*, 2023.
- [23] *Muhieddine Shebaro* and **Jelena Tešić**. Abcd: Algorithm for balanced component discovery in signed networks. In *Submitted to a Conference*.
- [24] *Muhieddine Shebaro* and **Jelena Tešić**. Scaling frustration index and corresponding balanced state discovery for real signed graphs. In *Submitted to a Conference*.
- [25] *Nicholas Warren*, *Ben Garrard*, *Elliot Staudt*, and **Jelena Tešić**. Transfer learning of deep neural networks for visual collaborative maritime asset identification. In *2018 IEEE 4th International Conference on Collaboration and Internet Computing (CIC)*, pages 246–255, Oct 2018.
- [26] Blake W Ford, Apan Qasem, **Jelena Tešić**, and Ziliang Zong. Migrating software from x86 to arm architecture: An instruction prediction approach. In *2021 IEEE International Conference on Networking, Architecture and Storage (NAS)*, pages 1–6, 2021.

- [27] Haitao Gong, **Jelena Tešić**, Jueqiang Tao, Xiaohua Luo, and Feng Wang. Automated pavement crack detection with deep learning methods: What are the main factors and how to improve the performance? *Transportation Research Record*, page 03611981231161358.
- [28] Andrew Magill, *Lia Nogueira de Moura*, *Maria Tomasso*, *Mirna Elizondo*, and **Jelena Tešić**. Enriching content analysis of tweets using community discovery graph analysis. In *Proceedings of the MediaEval 2020 Workshop*, volume 2882, 2020.
- [29] Taylor Mauldin, Anne H. Ngu, Vangelis Metsis, Marc E. Canby, and **Jelena Tešić**. Experimentation and analysis of ensemble deep learning in iot applications. *2019 VLDB DMAH*, 5(1):133–149, 2019.
- [30] Lucas Rusnak and **Jelena Tešić**. Characterizing attitudinal network graphs through frustration cloud. *Data Mining and Knowledge Discovery*, 6, November 2021.
- [31] Hanie Samimi, **Jelena Tešić**, and Anne Hee Hiong Ngu. Patient-centric data integration for improved diagnosis and risk prediction. In *Heterogeneous Data Management, Poly-stores, and Analytics for Healthcare*, pages 185–195, Cham, 2019. Springer International Publishing.
- [32] **J. Tešić**, D. Tamir, S. Neumann, N. Rische, and A. Kandel. Computing with words in maritime piracy and attack detection systems. In Dylan D. Schmorow and Cali M. Fidopiastis, editors, *Augmented Cognition. Human Cognition and Behavior*, pages 434–444, Cham, 2020. Springer International Publishing.