

# AI and data driven development

# Overview of on-going research topics and activities

Helena H. Olsson

Software Center Reporting workshop, Gothenburg, December 5th, 2019

## **Research Themes**



### Holistic DevOps Framework





Statistical comparison of black box algorithms for optimization of expensive functions

- Large scale comparison between algorithms
  - Over 25 algorithms
  - Over 170 benchmark functions
  - Different measurements metrics
  - Different conditions (noise, budget etc...)
- Applications
  - Field optimization
  - Hyperparameter tuning of ML models
  - AutoML

# Business transformations through ML experiments



## Data Management Challenges

#### MAPPING BETWEEN DATA MANAGEMENT CHALLENGES AND USE CASES

		Use cases of DL components					
Phase	Challenge	RS <sup>1</sup>	WPP <sup>2</sup>	HPP <sup>3</sup>	$MD^4$	FFD <sup>5</sup>	$MS^6$
Data Collection	Lack of metadata	Х	Х	Х	Х	Х	Х
	Data Granularity		Х	Х			
	Shortage of diverse samples		Х	Х	Х	Х	
	Need for sharing and tracking techniques	Х	Х	Х	Х	Х	
	Data Storage	Х					
Data Exploration	Statistical Understanding		Х			Х	Х
	Deduplication Complexity	Х	Х	Х	Х	Х	Х
	Heterogeneity in data	Х	Х	Х	Х	Х	
Data Preprocessing	Dirty data	Х	Х	Х	Х	Х	Х
	Managing sequences in data					Х	Х
	Managing categorical data				Х	Х	
Dataset Preparation	Data Dependency	Х	Х	Х	Х	Х	Х
	Data Quality	Х	Х	Х	Х	Х	Х
Data Testing	Tooling	Х	Х	Х	Х	Х	Х
	Expensive Testing	Х			Х		Х
Deployment	Data Extraction Methods	Х	Х	Х	Х	Х	Х
	Overfitting				Х	Х	
Post Deployment	Data sources and Distribution	Х	Х	Х			
	Data drifts	Х	Х	Х			
	Feedback loops	Х					

### Data pipeline robustness

- A robust data pipeline can be an effective solution to solve some of the data management problems
- **RQ:** How to keep a data pipeline healthy?
  - Automation of the pipeline
  - Monitoring of each step of the pipeline
  - Investigation of the potential errors that might occur at each step and implement detection mechanism
  - Implementation of automated alerting mechanisms and mitigation actions

### DataOps

- DevOps manages code, infrastructure and tools
- DataOps adds a fourth dimension data
  - DataOps automates the sequence of steps taken to deliver value to the customer where possible, minimize waste and redundancy, and stimulate continuous improvement
  - DataOps can support data teams evolve from environments of data silos, backlogs, and quality control issues to an agile, automated, and accelerated data cycle that continuously improves and delivers value to the customers

### Dynamic Data Management Challenges\*

#### Focus:

- Challenges involved in storing and preparing data for AI applications in embedded systems
- Collection, storing, pre-processing of data and training of AI applications in distributed architectures
- Data storage and analysis (edge server cloud)

#### **Challenges:**

- Expensive and error-prone collection of sensor data
- Difficulties in maintaining semantics of continuously evolving data
- Unclean and noisy data
- Restrictive security constraints
- Difficulties in interpreting heterogeneous and dynamic data sets
- Lack of well-defined purpose and goals

\*Ouhaichi H., Olsson H.H., Bosch J. (2019) Dynamic Data Management for Machine Learning in Embedded Systems: A Case Study. In: Hyrynsalmi S., Suoranta M., Nguyen-Duc A., Tyrväinen P., Abrahamsson P. (eds) Software Business. ICSOB 2019. Lecture Notes in Business Information Processing, vol 370. Springer, Cham

### Architectural distribution framework

### **Key factors:**

- Latency
- Privacy
- Convergence
- Accuracy
- Communcation cost
- Computational resource
- Data quality



### **Developing ML/DL Models: A Design Framework**

- Organizations integrate ML/DL technologies in software-intensive systems to increase value delivery
- Despite the level of expertise, developing high-performing and accurate ML/DL is a challenging task
- Need for a structured and systematic design process to build and operate ML/DL models
- Design process that identifies the phases involved in development of ML/DL models, the iterations that occur in between these phases and the challenges associated with each phase



### Phases, activities and challenges when developing ML/DL models

Phases	Activities	Challenges
Business case Specification	<ul> <li>Spend 1-2 days running simple models to check business case</li> <li>Discussion with Product owners</li> <li>Verify accessibility and availability of data</li> </ul>	<ul> <li>High setup costs</li> <li>Communication gap with Stakeholders</li> <li>High Al demands</li> <li>Data scientists shortage</li> <li>Dataset availability</li> </ul>
Data Exploration	<ul> <li>Meetings with domain experts</li> <li>Build and test hypothesis around data</li> <li>Using visualization techniques</li> </ul>	<ul> <li>Privacy concerns</li> <li>Noisy data</li> <li>Domain Experts accessibility</li> <li>Labelling</li> </ul>
Feature Engineering	<ul> <li>Start with high-dimensional features and scale down to few-dimensions</li> <li>Add relevant features scaling from small to large feature set</li> <li>Add features that are not directly part of the dataset</li> </ul>	<ul> <li>Increasing complexity</li> <li>Improper feature selection</li> </ul>
Experimentation	<ul> <li>State-of-the-art approach to find current best algorithm</li> <li>Obtain basic performance and experimenting with other algorithms</li> <li>Automate experimentation using tools</li> </ul>	<ul> <li>Introduce bias</li> <li>Selection of algorithms</li> <li>High DL complexity</li> <li>DL knowledge need</li> <li>Validity of related work</li> </ul>
Development	<ul> <li>Hyperparameter Tuning</li> <li>Final model selection depends on end goal requirements - Accuracy, Prediction time, Computational resources, Explainability, etc.</li> </ul>	<ul> <li>Difficulty in determining final model</li> <li>Model execution environment</li> <li>Complex hyperparameter settings</li> <li>Verification and Validation</li> </ul>
Deployment	<ul> <li>Place requirements on models before deployment</li> <li>Prepare code ready to put in docker container</li> <li>Plan for integration with internal systems</li> <li>Use API service reusable functionality to wrap the model</li> </ul>	<ul> <li>Less DL deployment</li> <li>Integration Problems</li> <li>Internal deployment</li> </ul>
Operational	<ul> <li>A/B testing</li> <li>Model Execution environment with built-in A/B support, Canary selection, etc.</li> <li>Services deployed have training interface, inference and evaluation interface</li> <li>Continuous retraining by adding functionality to existing models</li> </ul>	<ul> <li>End user interaction</li> <li>Model drifts</li> <li>Training-Serving Skew</li> </ul>

### **Artificial Intelligence**



## Challenges ML/DL Evolution

	Experiment Prototyping	Non-critical deployment	Critical deployment	Cascading deployment
assemble dataset	Issues with problem for- mulation and specifying de- sired outcome	Data silos, scarcity of la- belled data, imbalanced training set	Limitations in tech- niques for gather- ing training data from large-scale, non-stationary data streams	Complex and effects of data dependencies
create model	Use of non- representative dataset, data drifts	No critical analysis of training data	Difficulties in build- ing highly scalable ML/DL pipeline	Entanglements causing difficul- ties in isolating improvements
train and evaluate model	Lack of well- established ground truth	No evaluation of models with business-centric measures	Difficulties in repro- ducing models, results and debugging DL models	Need of tech- niques for sliced analysis in final model
deploy model	No deployment mechanism	Training- serving skew	Adhering to strin- gent serving require- ments e.g., of latency, throughput	Hidden feedback- loops and unde- clared consumers of the models

Table 2. Challenges in the evolution of use of ML/DL components in software-intensive systems

### **Selected** publications

Olsson, H.H., and Bosch, J. (2019). Data Driven Development Adoption Process. In Proceedings of the 20th International Conference on Product-Focused Software Process Improvement (PROFES), November 27 – 29, Barcelona, Catalunya, Spain.

Raj, A., Bosch, J., Olsson, H.H., Arpteg, A. and Brinne, B. (2019). Data Management Challenges for Deep Learning. In Proceedings of the EUROMICRO Conference on Software Engineering for Advanced Applications, August 28 - 30, Kallithea, Chalkidiki, Greece.

Lwakatare, L. E., Raj, A., Bosch, J., Olsson, H. H., & Crnkovic, I. (2019). A Taxonomy of Software Engineering Challenges for Machine Learning Systems: An Empirical Investigation. *In International Conference on Agile Software Development (pp. 227-243). Springer, Cham.* 

Mattos, D.I., Bosch, J., Olsson, H.H., Dakkak, A., and Bergh, K. (2019). Automated Optimization of Software Parameters in a Long-Term Evolution Radio Base Station. *In Proceedings of the IEEE Annual Systems Conference, April 8 – 11, Orlando, Florida, USA*.

Mattos, D.I., Bosch, J., and Olsson, H.H (2019). Multi-armed bandits in the Wild: Common Pitfalls in Online Experiments. *Information and Software Technology*, pp.1-14.

Mattos, D.I., Bosch, J., and Olsson, H.H. (2019). ACE : Easy Deployment of Field Optimization Experiments. *In the Proceedings of the European Conference on Software Architecture (ECSA), September 9 – 13, Paris, France.* 

Mattos, D.I., Bosch, J., and Olsson, H.H. (2019). Leveraging Business Transformations with ML Experiments. *In Proceedings of the International Conference on Software Business (ICSOB), November 18 – 20, Jyväskylä, Finland.* 

Figalist I., Elsner C., Bosch J., Olsson H.H. (2019) Scaling Agile Beyond Organizational Boundaries: Coordination Challenges in Software Ecosystems. In: Agile Processes in Software Engineering and Extreme Programming. XP 2019. Lecture Notes in Business Information Processing, vol 355. Springer, Cham.

Figalist I., Elsner C., Bosch J., Olsson H.H. (2019) Business as Unusual: A Model for Continuous Real-Time Business Insights Based on Low Level Metrics. *In: Proceedings of the 45th EUROMICRO Conference on Software Engineering and Advanced Applications, Kallithea, Greece.* 

Figalist I., Elsner C., Bosch J., Olsson H.H. (2019) Customer Churn Prediction in B2B Contexts. In Proceedings of the International Conference on Software Business (ICSOB), November 18 – 20, Jyväskylä, Finland.

Ouhaichi, H., Olsson, H.H., and Bosch, J. (2019). Dynamic Data Management for Machine Learning in Embedded Systems: A Case Study. In Proceedings of the International Conference on Software Business (ICSOB), November 18 – 20, Jyväskylä, Finland.