

## On Learning Data Representation

Guillermo D. Canas, [guilledc@MIT.EDU](mailto:guilledc@MIT.EDU)

Istituto Italiano di Tecnologia and Massachusetts Institute of Technology

Lorenzo Rosasco\*, [lrosasco@MIT.EDU](mailto:lrosasco@MIT.EDU)

University of Genova, Istituto Italiano di Tecnologia and Massachusetts  
Institute of Technology

Machine learning has played a central role in the recent successes of artificial intelligence and provided fundamental tools for analyzing complex, high dimensional data in Science and Engineering. A main bottleneck of current machine learning systems is the large amount of labeled data required, which in turn calls for intense computations and considerable human supervision. Data representation is broadly acknowledged to be key in solving the above problems. The intuition that these algorithms adaptively reduce the degrees of freedom in the data, and lead to improved performances, is confirmed by promising empirical results. However, these intuitions and empirical observations are not grounded in solid theoretical foundations, a fact that arguably hinders progress in the field. In this work, we discuss an approach to establish a statistical theory of learning representations from data akin to the classic supervised statistical learning theory. The theory we propose borrows ideas from Signal Processing, Machine Learning and Statistics, while establishing novel bridges between such diverse field as Optimal Quantization, Optimal Transport Theory. Computational solutions for learning representations from data, based on piecewise models are studied.

Keywords: Statistical Learning Theory, Data Representation, Unsupervised Learning.