

Campus Luzern – A Cooperative Project of the University-Level Institutions in Lucerne

## «*Big Data Analytics*»

<b>Tutor</b>	Luigi Curini, University of Milan
<b>Organization</b>	Campus Luzern
<b>Language</b>	English
<b>ECTS-Points</b>	4
<b>Content</b>	<p>Big data are those labeled, for strange reasons, with the capitalized “Big”. Nevertheless, they are still “data”, although with some specific characteristics: large volume, high frequency and, most notably, unpredictability - data come in the many different forms, they are raw, messy, unstructured, not ready for processing, and so on. Still, these data convey a lot of information to social scientists and good statistical techniques are required in order to extract meaningful results from them. In this workshop we will focus on a specific type of Big data, namely digital texts, both from social media as well as other sources (such as legislative speeches or electoral programs). The aim is to provide an introductory guide to this exciting new area of research, while also offering guidelines on how to effectively use statistical methods on texts for social scientific research by discussing the advantages, but also the limits, of each approach. The attention will be devoted to three main areas: 1) scaling methods that allow to estimate the location of actors in some policy space; 2) supervised classification methods, including machine learning algorithms, that allow to organize texts into a set of pre-defined categories; 3) unsupervised classification that allow to discover new ways of organizing texts into a set of unknown categories. Time permitting, beyond the Bag-of-Word approach we will also briefly cover the word-embedding approach.</p> <p>Day 1: supervised and unsupervised scaling methods (Wordscores and Wordfish)</p> <p>Day 2: unsupervised classification methods (LDA and Structural Topic Models)</p>

	<p>Day 3: an introduction to supervised classification methods &amp; on how to extract texts from social media source</p> <p>Day 4: supervised classification methods (machine learning algorithms)</p>
<p><b>Prerequisites/ Materials (opt.)</b></p>	<p>An elementary knowledge of R (having attended any of the introductory workshops offered by Campus Lucerne usually satisfies this requirement), plus a curiosity towards applied statistics, are good prerequisites for the lab sessions. Participants will familiarize with <code>quanteda</code>, one of the most well-known and better-developed text-mining R package. On top of that, in our lab examples we will employ several other packages, in particular when discussing about classification methods (for example: <code>stm</code>, <code>topicmodels</code>, <code>naivebayes</code>, <code>e1071</code>, <code>randomForest</code>, <code>xgboost</code>). We will also use some R packages to extract texts from social media source, such as <code>rtweet</code> and <code>tuber</code>. All the datasets, replication files of the lab sessions and reference texts will be made available at a dedicated URL before the beginning of the workshop. Workshop participants should bring their own laptop with R, RStudio and the relevant packages previously installed and functioning (instructions will be circulated beforehand).</p>
<p><b>Charge</b></p>	<p>This Campus Luzern offering is directed at researchers, post-docs, doctoral students, people interested in starting a dissertation and employees of Lucerne universities and is free of charge for these persons.</p> <p>Certificates with the respective ECTS-Points are prepared on request (<a href="mailto:graduate.resources@campus-luzern.ch">graduate.resources@campus-luzern.ch</a>).</p>