

## ADHOC-TE-2016-0892 STATE of the ART - DIGITAL DATA COLLECTION

WEB PAGES	DATA SOURCES
<p>On the World Wide Web there are static and dynamic websites</p> <p><b>Static Websites</b> are characteristic for the beginning of the Web and consist of webpages (HTML files) hosted on Web servers and accessible through a browser.</p> <p>The content and form (mark-up elements) are not separate. The structure of the pages is fixed and any change of the structure, content or form is done through editing the individual files or page templates (in some cases) or stylesheets. Due to better loading speeds, static websites have seen a resurgence in the past few years through many user-friendly site-makers / generators.</p> <p><b>Dynamic Websites</b> are generated by applications running on a server and connect to a database (hosted on the same server or another) where the content and configuration/structure are stored.</p> <p>These websites are created and managed through WEB CONTENT MANAGEMENT SYSTEMS.</p> <p>Individual pages that are accessed through a browser do not exist on a Web server, but are generated for each request using templates which combine HTML, scripting languages (PHP, JavaScript) and filling them in with content from the database, thus generating each page dynamically, adapted to the context.</p>	<p>OPEN DATA (datasets made public or publicly accessible and useable without restrictions)</p> <p>Open data are data which can be used freely, reused, redistributed by anyone – subjected at the most to attribution</p> <p>DATA FROM DOCUMENTS (usually, tables in PDF)</p> <p>DATA FROM WEBSITES</p> <p>DATA FROM API INTERROGATION (Facebook, Twitter, Wikipedia etc.)</p> <p>API (Application Programming Interface) – sets of functions which are made available for developers who want to integrate a system with another. To encourage the development of Web or mobile apps, many social media systems have (partially) open APIs</p> <p>According to the type of data and their specific format, several automated data collection tactics may be used to gather large datasets</p> <ul style="list-style-type: none"> <li>• OCR (Optical Character Recognition) and document conversion</li> <li>• Web Scraping, Web Wrapping and Web Crawling</li> <li>• API Interrogation</li> </ul>
WEB SCRAPING	API INTERROGATION
<p>Web scraping / Screen scraping refers to different methods that can be used to collect data from Web pages, usually from dynamic web pages. It is also sometimes called Web data extraction, screen scraping or Web harvesting.</p> <p>It is a form of DATA MINING and it is based on the identification of distinct patterns in HTML and CSS templates of dynamic pages to detect information of the same type and to automate the data collection process</p> <p>The stages of a scraping project:</p> <ul style="list-style-type: none"> <li>• <b>Choosing the target content</b></li> <li>• <b>Defining selectors / patterns</b></li> <li>• <b>Simulating navigation</b></li> <li>• <b>Configuring automation</b></li> </ul> <p><b>Tools:</b> Chrome Web Scraper, Outwit Hub, Octoparse, Helium Scraper, Import.io, Screen Scraper</p> <p><b>Choosing the target content</b></p> <p>Search results, e-commerce websites, media sharing platforms, archives of news websites or blogs have linear structures and are usually displayed in a paged list (with navigation – `next` or `show more` buttons)</p> <p>Forums or threaded comments have tree-like structures.</p> <p>In dynamic websites, the same HTML tags and CSS styles are used for content of the same type.</p> <p><b>Defining selectors / patterns</b></p> <p>The semi-structured nature of webpages may be used.</p> <p>HTML pages may be represented as tree structures of nested HTML element tags</p> <p>Each level of the hierarchy represents a nesting level of the HTML elements. This representation is called DOM (Document Object Model)</p> <p>For each HTML element in a document, a unique path can be defined in a similar way in which we call up the path of a file in a filesystem.</p> <p>For example:</p> <p>/html/head/title will refer to the title of the document</p> <p>/html/body/div[2] will refer to the second div element in the body</p> <p>Xpath specifications allow us to refer to the surrounding elements of a given element, conceptualized as a node with certain relations (ancestor, parent, sibling, child, descendant, preceding, following) in the document tree structure.</p> <p><b>Types of navigation</b></p> <ul style="list-style-type: none"> <li>• Linear structures – Pagination (blogs/news)</li> <li>• Tree structures – Multilevel (forums/comments)</li> <li>• Network structures – Graph (wiki type sites)</li> <li>• Tabular structures – sortable by variables (e-commerce)</li> <li>• Relational structures – entities of several types (IMDB)</li> </ul>	<p>API = Application Programming Interface is a software-to-software interface that allows data exchange between to applications</p> <p>For example, Facebook offers the public API OpenGraph so that other applications created by developers to be able to integrate functionality/personalized information from Facebook – Facebook authentication, friends who liked a certain content etc.</p> <p><b>Types of API</b></p> <p>RESTful- (Representational State Transfer) – the most used form of API – communication through HTTP</p> <p>SOAP (Simple Object Access Protocol) – data transfer in XML format</p> <p><b>Types of HTTP requests</b></p> <p>GET – to get data</p> <p>PUT – to edit existing data</p> <p>POST – to add new data</p> <p>DELETE – to delete data</p> <p>Most APIs require authentication.</p> <p>HTTP Basic Access Authentication – user and password are transmitted in the header of the HTTP request</p> <p>OAuth 1.0/2.0 – a unique token is generated for the user</p> <p>A significant number of APIs will transmit responses as JSON (JavaScript Object Notation) - the most widely used data format for data interchange on the web</p> <p><b>Facebook Graph API</b></p> <p>Nodes – users, pages, posts, comments</p> <p>Edges – connections between nodes – the comments of a certain post, the posts of a certain page</p> <p>Fields – information about entities, objects – pages, users, posts</p> <p>Each node has a unique Object ID</p> <ul style="list-style-type: none"> <li>• Usually APIs limit the quantity of data that can be transmitted through a single response</li> <li>• Paginated responses may allow the collection of large quantities of data</li> <li>• In the aftermath of the Cambridge Analytica scandal, Facebook has introduced severe limitations in the access of data through its API (including data from public Facebook pages)</li> </ul> <p><b>Tools: Facepager</b> – usage outline:</p> <p>Creating a database</p> <p>Configuring an interrogation preset</p> <p>Configuring data columns</p> <p>Authentication and token</p> <p>Starting the interrogation (actual data collection)</p> <p>Exporting the data in CSV</p> <p>Scenarios: Facebook, Twitter, YouTube, Wikipedia, generic APIs</p> <p><b>Alternative tools:</b> API interrogation with cURL command line tool, python libraries for interacting with APIs etc.</p>

## ADHOC-TE-2016-0892 STATE of the ART DIGITAL DATA ANALYSIS

DATA CLEANUP AND CONVERSION	ELEMENTS OF COMPUTATIONAL LINGUISTICS
<p>Data collection from online sources using Web scrapers or API interrogation will usually mean data is initially in text format, character strings.</p> <p>Data extracted through Web Scraping or API interrogation are semi-structured data.</p> <p>Issues with saving or exporting data:</p> <p>File format</p> <p>Maximum field sizes (for each tool used)</p> <p>Character encoding (special characters / diacritical marks may not be recognized by some analytical software)</p> <p><b>Data types</b></p> <ul style="list-style-type: none"> <li>• character</li> <li>• string</li> <li>• integer</li> <li>• float</li> <li>• boolean (TRUE/FALSE)</li> <li>• Date / Time</li> <li>• Coordinates (latitude/longitude)</li> </ul> <p><b>File formats</b></p> <p>Character Separated Values (comma, tab, semicolon)</p> <p>Fixed width – each column takes up a number of characters on each row</p> <p>Importing CSV files will require defining the field delimiter and any other characters used to define data fields</p> <p>When importing text files into different analysis tools, after defining the data columns, fields are sometimes associated a certain data type. When converting data, maximal values that can be stored into each data type should be taken into consideration.</p> <p>For data resulting from scraping several different websites (for example news sites) the date/time format will often be different and will need to be converted to a common format.</p> <p><b>Some common string operations</b></p> <p>CLEAN – deleting unprintable characters</p> <p>TRIM – deleting extra spaces from the beginning or the end</p> <p>UPPER/LOWER – conversion to upper or lower case</p> <p>LEN- length of a character string</p> <p>CONCAT / CONCATENATE / TEXTJOIN – joins two or more strings</p> <p>LEFT / RIGHT – a certain number of characters from the beginning or end of a string</p> <p>REPLACE / SUBSTITUTE – replaces a string with another</p> <p>CONTAINS – returns TRUE if a string contains a given substring</p> <p><b>Tools: Tableau Data Prep</b> – large dataset overview, filtering, sorting, pivoting, character/string replacing, calculated fields etc.</p> <p><b>ASAP Utilities</b> – replacing accented characters, deleting extra spaces, recognizing dates, changing date/time formats</p>	<p><b>Stopwords</b> – words in a language that are usually excluded before working with natural language processing tools – usually very frequent words or particles, deictics, pronouns, conjunctions, compound verb forms</p> <p><b>POS tagging</b> – Identification of parts of speech</p> <p><b>Stemming</b> – process through which the endings of words are trimmed in hopes of reducing words to their base form.</p> <p><b>Lemmatization</b> – process that uses a vocabulary and morphological analysis to identify and reduce all forms of a word to their base forms.</p> <p><b>Jaccard Distance/ Jaccard similarity coefficient</b> is used to compare elements of two sets with the purpose of observing which are common and which are distinct. It is a measure of similarity of two data sets or codes</p> <p><b>Hierarchical clustering</b> - data mining methods that group data into categories and subcategories according to a measure of distance / difference between the sets and their elements.</p> <p>Co-occurrence analysis – a co-occurrence matrix describes in a binary way (presence/absence) an object with respect to a context (a word or a code in a document/text)</p> <p><b>Co-occurrence network</b> – a more visually intuitive representation based on a co-occurrence matrix where nodes represent objects (words/codes), and edges/lines between them represent their presence in the same context (co-occurrence). Nodes or edges may be weighed (sized or styled with according to their frequency or relative distance score)</p> <p><b>TF</b> - term frequency - (word or code)</p> $TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$ <p><b>IDF</b> - inverse document frequency - - a measure of the importance of a term. While computing TF, all terms are considered equally important. However, it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:</p> $IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$ <p><b>Tools:</b></p> <p><b>KH Coder</b> – visual interface for natural language processing (based on R) for co-occurrence networks, hierarchical clustering of words or codes, heatmaps etc.</p> <p><b>Sketch Engine</b> – corpus linguistics, lexicography</p> <p>Spacyr and Shiny – R libraries for visualizing co-occurrence networks based on spacy tokenization</p>
<p><b>KNOWLEDGE REPRESENTATION AND INFORMATION RETRIEVAL</b></p> <p><b>Semantic networks</b> are a logic-based formalism for knowledge representation. Semantic networks are graphs which are constructed from both a set of vertices (or nodes) and a set of directed and labeled edges. The vertices or nodes represent concepts, and the edges represent semantic relations between the concepts. Knowledge about accepted meanings should be processed in adjacent regions of the semantic network. Therefore, semantic networks are often termed “associative networks.”</p> <p><b>Formal Concept Analysis (FCA)</b> is a method of knowledge representation introduced in the 1980s by Rudolf Wille, rooted in the pragmatic philosophy of Charles Sanders Peirce, based on a binary incidence relation, and building on applied lattice and order theory. It has applications in various fields and its advantage lies in the possibility to visualize and explore formal concepts in a formal context (a data table that represents binary relations between items in a set of objects and items in a set of attributes) as representations of complete lattices.</p> <p><b>Tools:</b> FCA Tools Bundle</p>	<p><b>MACHINE LEARNING and NATURAL LANGUAGE PROCESSING</b></p> <p><b>Word embedding</b> is one of the most popular representation of document vocabulary. It is capable of capturing context of a word in a document, semantic and syntactic similarity, relation with other words, etc.</p> <p>The objective is to have words with similar context occupy close spatial positions. Mathematically, the cosine of the angle between such vectors should be close to 1, i.e. angle close to 0.</p> <p><b>Word2Vec</b> is a two-layer neural network trained to reconstruct linguistic contexts of words. It is a method to construct such an embedding. It uses large corpus inputs and produces a vector space of several hundred dimensions, each word being assigned a vector in the space. The word vectors' positioning so that words which have common contexts in the corpus are located in proximity to one another.</p> <p><b>NER – Named Entity Recognition, Text Categorization/Classification</b> – require training language models or using pre-trained models.</p> <p><b>Tools:</b> Gensim (Word2Vec), Spacy (NLP), Prodigy (Annotation interface), PyTorch, LASER, scikit-learn</p>