

Analyse Visuelle de Réseaux Sociaux Historiques: Traçabilité, Exploration et Analyse

*Visual Analytics for Historical Social Networks:
Traceability, Exploration, and Analysis*

**Thèse de doctorat de l'université Paris-Saclay et de
Telecom Paris**

École doctorale n°580 : Sciences et technologies de l'information et de la
communication (STIC)

Spécialité de doctorat: Informatique

Graduate School : Informatique et Sciences du Numérique

Référent : Faculté des sciences d'Orsay

Thèse préparée au Laboratoire interdisciplinaire des sciences du numérique
(Université Paris-Saclay, CNRS, Inria), et à Telecom Paris, sous la direction de
Jean-Daniel FEKETE, Directeur de recherche et la co-direction de Christophe
Prieur, Professeur des universités.

Thèse soutenue à Paris-Saclay, le xx décembre 2022, par

Alexis PISTER

Composition du jury

Ulrik Brandes

Professeur, ETH Zürich

Guy Melançon

Professeur, Université de Bordeaux

Wendy Mackay

Directrice de recherche, Univ. Paris-Saclay, CNRS, Inria, LISN

Uta Hinrichs

Professeur, University of Edinburgh

Laurent Beauguitte

Chargé de recherche, CNRS

Jean-Daniel Fekete

Directeur de recherche, Univ. Paris-Saclay, CNRS, Inria, LISN

Christophe Prieur

Professeur, Université Gustave Eiffel

Rapporteur & Examineur

Rapporteur & Examineur

Examinatrice

Examinatrice

Examineur

Directeur de thèse

Directeur de thèse

Titre: Analyse Visuelle de Réseaux Sociaux Historiques: Traçabilité, Exploration et Analyse

Mots clés: 3 à 6 mots clefs (version en français)

Résumé:

Title: Visual Analytics for Historical Social Networks: Traceability, Exploration, and Analysis

Keywords: visual analytics, social network analysis, social network visualization, social history.

Abstract:

This thesis aim at identifying how Visual Analytics can support historians in their social network analysis process, from the collection of historical documents to the formulation of high-level socio-historical conclusions. Historical Social Network Analysis is a method followed by historians to study social relationships and interactions between groups of actors (families, institutions, political elites, companies, etc.) to understand their underlying structure while characterizing specific behaviours. Social historians are able to reconstruct relationships of the past using the rich information contained in historical document, such as marriage acts, migration forms, brith certificates, and census. Through visualization and analytical methods, they can describe the global structure of studied groups and explain individual behaviors through local network patterns. However, the inspection, encoding, correction, and modeling process of the historical documents leading to a finalized network is intricate and often results in inconsistencies, errors, distortions, simplifications, and traceability issues. Moreover, usability and analytical interpretation issues often limit the usage of visual interfaces in History. For these reasons, social historians are not always able to make thorough historical conclusions with current analytical and visualization tools. In this thesis, I aim to identify how visual analytics—the integration of data mining capabilities into visual interfaces with direct manipulation and interaction—can support social historians in their process, from the collection of their data to the answer to high-level historical questions. Towards this goal, I first formalize the workflow of historical network analysis in collaboration with social historians, from the acquisition of their sources to their final visual analysis and point that visual analytics tools supporting this process should satisfy traceability, simplicity, and document reality principles to ease bask and forth between the different steps, provide tools easy to manipulate, and not distort the content of sources with modifications and simplifications. Particularly, I propose to model historical sources into bipartite multivariate dynamic social networks with roles to satisfy those properties. This modeling allows a concrete representation of historical documents, hence letting users encode, correct, and analyze their data with the same abstraction and tools. Leveraging this data model, I propose two interactive visual interfaces to manipulate, explore, and analyze this type of data with a focus on usability for social historians. First, I present ComBiNet, which allows an interactive exploration leveraging the structure, time, localization, and attributes of the data model with the help of coordinated views, a visual query system, and comparison mechanisms. Finding specific patterns easily and comparing them, social historians are able to find inconsistencies in their annotations and answer their high-level questions. The second system, PK-Clustering, is a concrete proposition to increase the usability and effectiveness of clustering mechanisms in social network visual analytics systems. It consists in a mixed-initiative clustering interface that let social scientists create meaningful clusters with the help of their prior knowledge, algorithmic consensus, and interactive exploration of the network. Both systems have been designed with continuous feedback from social historians, and aim to increase the traceability, simplicity, and document reality of visual analytics supported historical social network research. I conclude with discussions on the potential merging of both systems and more globally on research directions towards better integration of visual analytics systems on the whole workflow of social historians. Such systems with a focus on those properties—traceability, simplicity, and document

reality—can limit the introduction of bias while lowering the requirements for the use of quantitative methods for historians and social scientists which has always been a controversial discussion among practitioners.

Contents

1	Introduction	13
1.1	Social History and Historical Social Network Analysis	14
1.2	Visualization and Visual Analytics	16
1.3	Visual Analytics Supported Historical Network Research	19
1.4	Contributions and Research Statement	21
2	Related Work	23
2.1	Notations	24
2.2	Visualization	26
2.2.1	Information Visualization	26
2.2.2	Visual Analytics	29
2.3	Quantitative Social History	29
2.3.1	History, Social History, and Methodology	30
2.3.2	Quantitative History	31
2.3.3	Digital Humanities	32
2.4	Historical Social Network Analysis	34
2.4.1	Sociometry to SNA	35
2.4.2	Methods and Measures	36
2.4.3	Historical Social Network Analysis	38
2.4.4	Network Modeling	40
2.5	Social Network Visualization	41
2.5.1	Graph Drawing	41
2.5.2	Social Network Visual Analytics	43
3	Historical Social Network Process, Pitfalls, and Network Modeling	45
3.1	Context	46
3.2	Related Work	47
3.2.1	History Methodology	47
3.2.2	Historian Workflows	48
3.3	Historical Social Network Analysis Workflow	48
3.3.1	Examples	49
3.3.2	Workflow	50
3.3.3	Visual Analytics Supported Historical Social Network Analysis	52
3.4	Network Modeling and Analysis	54
3.4.1	Network Models	54
3.4.2	Bipartite Multivariate Dynamic Social Network	57
3.5	Applications	59
3.6	Discussion	59
3.7	Conclusion	61

4	ComBiNet: Visual Query and Comparison of Bipartite Dynamic Multivariate Networks with Roles	63
4.1	Context	64
4.2	Related Work	66
4.2.1	Graphlet Analysis	66
4.2.2	Visual Graph Querying	67
4.2.3	Visual Graph Comparison	67
4.2.4	Provenance	67
4.3	Task Analysis and Design Process	68
4.3.1	Use Cases	68
4.3.2	Tasks Analysis	70
4.4	The ComBiNet System	72
4.4.1	Visualizations	73
4.4.2	Query Panel	74
4.4.3	Comparison	81
4.4.4	Implementation	83
4.5	Use Cases	84
4.5.1	Construction sites in Piedmont (#1)	84
4.5.2	French Genealogy (#2)	84
4.5.3	Marriage acts in Buenos Aires (#3)	87
4.5.4	Sociology thesis in France	88
4.6	Formative Usability Study	90
4.6.1	Feedback	91
4.7	Discussion	92
4.8	Conclusion and Future Work	93

List of Figures

1.1	Business contract originated from Nantes (France) during the 17th century. See [34] for more detail of the historian process to analyze her sources.	16
1.2	Marriage, partnership, trading, banking, and real estate networks of the powerful families of Florence from [105]. We can see the central position in the network of the Medici Family.	17
1.3	Abstraction of the VA process. It is characterized by continuous interactions between the data, visualizations, models, and knowledge. Image from [71]. . . .	18
1.4	Node-link diagram of a medieval social network of peasants, produced with a force-directed layout, commonly used in SNA softwares. Image from [?].	20
2.1	Categorization of visual variables which can be used to represent network data, resulting in many different network representations. Image from [8].	27
2.2	Anscombe quartet. The four datasets have the same descriptive statistics (average, variance, correlation coefficient) but very different structures. Image from [4].	28
2.3	TULIP software designed for application-independant network visual analytics [?]. The view shows a dataset among multiple interactive coordinated views. Users can also apply data mining algorithms on the data to extract interesting patterns.	30
2.4	Correspondence letters of Benjamin Franklin and his close relationships, using a map and an histogram, accessible online on the republic of letter website [?]. . .	34
2.5	Moreno's original sociogram of a class of first grades from [94] (left). The diagram shows 21 boys (triangles) and 14 girls (circles). The same sociogram using modern practices generated from Gephi from [53] (right). The color encodes the number of incoming connections.	36
2.6	All possible graphlets of size 2 to 5 for undirected graphs	37
2.7	Cicero personal communication network represented with a node-link diagram. Image from [?]	39
2.8	Different criteria are proposed to enhance node-link diagram readability. Image from [75]	42
2.9	NodeTrix system showing a scientific collaboration social network with clusters. Each cluster is represented as a matrix, Image from [64].	43
2.10	Vistorian interface [128] used to explore a historical social network of business trades in the 17th century, with a coordinated node-link diagram and a matrice view.	44
3.1	HSNA workflow is split into five steps: textual sources acquisition, digitization, annotation, network creation, and network visualization/analysis. Practitioners typically have to do back and forth during the process. I list potential pitfalls for each step.	50

3.2	Three properties essential to VA systems supporting the social historians workflow: <i>traceability</i> , <i>document reality</i> , and <i>simplicity</i>	53
3.3	bipartite multivariate dynamic network modeling for two cases of marriage acts of example #3. Some marriage acts mention the parents of the spouses, which is a relationship different than the marriage in itself. This case can be modeled using a document model (a) or an event model (c) by splitting the document into several different event nodes. The other case refers to documents that do not mention the parents (b) and in that case, the network represents both the documents and the events with the same model. M: Marriage, H: Husband, W: Wife, T: Witness, (H/W)(M/F): Husband/Wife Mother/Father. Yellow links refer to parenting mentions/relationships.	60
4.1	The ComBiNet system used to compare two subgroups of a social network of contracts from [25], extracted with dynamic visual queries. (A) and (B) show the two visual queries created by the user in the query panel using an interactive node-link diagram editor (V6), dynamic query widgets (V7), and the equivalent Cypher script (V8). The right part shows ComBiNet's global interface in <i>comparison</i> mode: (V1) Network visualization panel, (V2) Map of the geolocalized nodes, (V3) Table of persons, (V4) Graph measures comparison, (V5) Attribute distribution plots, and (V9) Provenance tree. The two visual queries on the left, translated into Cypher queries below, select the "Menafoglio" family on the left, and the "Zo" family on the right, along with their construction contracts and close collaborators.	72
4.2	ComBiNet interface wreal-timeith the dataset of collaboration #1. The user selected the <code>__year</code> attribute, showing the distribution of document years with a histogram (bottom right), and coloring the documents node on the bipartite view (left) and map view simultaneously (right).	75
4.3	All link creation possibilities: Any link type (left), one selected link type, here guarantor (middle left), the union of several link types (middle right), several links with different types (right)	76
4.4	Visual queries created to answer questions 2 and 6 of our collaboration #1. (a) The visual query retrieves individuals who are mutually guarantors to each other in separate construction contracts. (b) The two visual queries retrieve the documents—along with the signatories—of Torino (<i>Turin</i> in french) (left) and of Torino surroundings (<i>Turin Territoire</i> and <i>Piemont</i>) (right)	77
4.5	Widget designs for the different attribute types: checkboxes for categorical attributes (top), text input for nominal attributes (middle), and a double slider for numerical attributes (bottom). The categorical attribute example shows the inputs letting users create new constraints for other attributes and other nodes.	78

4.6	Results of question 2 of collaboration #1: (a) shows a subset of the table view with every occurrence of the pattern found. (b) shows the summary panel, with the graph measures and the attributes view with the <i>origin</i> attribute selected and the Sankey option checked. It allows us to see the attribute distribution of the persons included in the pattern and see if there is a relationship between persons who are mutually guarantors and their origin.	79
4.7	Two ways of showing the distribution of “type chantier” (type of works), a categorical attribute with three possible values “ <i>religious</i> ”, “ <i>military</i> ”, and “ <i>civilian</i> ”. (a) A query matching the contracts made by the same person (<i>per1</i>) as an “approbator” (green link to <i>doc2</i>) after being a “guarantor” (blue link to <i>doc1</i>) using the constraint (<i>doc2._year</i> > <i>doc1._year</i>). (b) Stacked bar chart for the matches, the earlier contract (<i>doc1</i>), the older contract (<i>doc2</i>), and (c) Sankey diagram with the early values on the left and the last on the right. The Sankey diagram reveals the value changes between the two documents: the guarantor who worked initially on religious work switched to military work. . . .	80
4.8	Provenance tree to answer question 2 of collaboration #1: left branch leads to Torino documents (the node is labeled as A) while right branch leads to surrounding documents (the node is labeled as B). The user hovers over one node, revealing a tooltip that shows the visualization of the node’s query.. . . .	81
4.9	Comparison table of the network measures for Torino subgraph (A) and Torino surroundings subgraph (B).	82
4.10	Distribution of the type of constructions, the years, and the betweenness centrality for the documents and signatories of Torino (A), Torino surroundings (B), and the whole graph (top).	83
4.11	Menafoglio (a) and Zo (b) families were retrieved with queries and highlighted in the bipartite node-link and map views.	85
4.12	Attributes distributions plots between the whole graph, the <i>Menafoglio</i> family (A), the <i>Zo</i> family (B), and $A \cap B$, for the <i>region</i> , <i>type_chantier</i> , <i>material type</i>	86
4.13	Map of the migrations in France which occurred across several generations.	87
4.14	Migrations across departments over three generations	87
4.15	Sankey diagrams showing the migration of people in the 18th and 19th centuries, extracted from their birth and death places.	88
4.16	ComBiNet used to request persons appearing as husband, wife, or witness in two marriages that occurred 70 years apart or more.	89
4.17	ComBiNet used for exploring theses of sociology defended in France between 2016 and 2021. The bipartite and map views show an overview of two visions of the network. The user selects the <i>region</i> attribute, showing the geographical distribution of the defended thesis.	90

4.18 Sociology thesis dataset explored with ComBiNet. The user constructed a visual query to see if there are symmetrical relationships between thesis directors and reviewers (or jury directors). The <i>region</i> attribute is selected with the Sankey option, letting the user see if there are correlations between the regions of the thesis found in this pattern.	91
---	----

List of Tables

2.1	Comparison table of most widely used visualization and analytical tool for HSNA. Visualizations: number of different visualization techniques, layout, and interactions. SNA and Models: Number of proposed SNA measures and algorithms. Clustering: Number of proposed clustering algorithms. Filtering: Possibilities of filtering according various criteria. Interaction/Direct Manipulation: Number of possible interactions mechanisms directly applicable on the visualizations. . . .	43
3.1	Resulting networks using different models produced by one document of the examples detailed in §3.3.1: co-occurrence, unipartite and bipartite models. The first column shows the partial transcription of real documents (simplification for collaboration #1). Colors represent annotations concerning the persons mentioned, their roles, and attributes. Underlines refer to information related to the events and which can be encoded as document/event attributes. Only the time is represented for simplification, but other attributes would follow the same schema. H: Husband, W: wife, T: Witness, M: Marriage, A_N : Associate, G: Guarantor, Ap: Approbator, C: Construction, F: Father, M: Mother, C: Child.	56
4.1	Tasks to support during exploration, according to our expert collaborators, are split into 3 main high-level tasks.	71
4.2	Comparison of the data model of several VA systems aimed at exploring bipartite social networks.	72

2 Related Work

Contents

2.1	Notations	24
2.2	Visualization	26
2.2.1	Information Visualization	26
2.2.2	Visual Analytics	29
2.3	Quantitative Social History	29
2.3.1	History, Social History, and Methodology	30
2.3.2	Quantitative History	31
2.3.3	Digital Humanities	32
2.4	Historical Social Network Analysis	34
2.4.1	Sociometry to SNA	35
2.4.2	Methods and Measures	36
2.4.3	Historical Social Network Analysis	38
2.4.4	Network Modeling	40
2.5	Social Network Visualization	41
2.5.1	Graph Drawing	41
2.5.2	Social Network Visual Analytics	43

Social historians rely on textual historical documents to study social groups through their structures and socio-economic characteristics in societies of the past [?]. They read and analyze documents they can find from a period and subject of interest, and make their conclusions through deep inspection and cross-referencing of the information they found. Several methods have been developed in History to extract and analyze the information contained in the documents in a methodical way [?], based on qualitative or quantitative methods—among which HSNA is now widely popular [?]. HSNA is a method consisting in modeling the relational information mentioned in the documents—such as family, business, or friendship ties—in a network, to be able to characterize and explain social behaviors through the description of the network’s structure [72, 149]. This approach is directly inspired by SNA, which is a well-known method that sociologist theorized to understand and describe real world social relationships modeled as networks [43, 127]. Historians appropriated this method, by extracting relationships from historical documents. The specifics of HSNA in contrast of its sociology counterpart is therefore the modeling of the network from the historical documents—which are at the core of the historical work [113]—and the integration of the temporal dimension which is often disregarded in traditional SNA but central in history. Once they successfully constructed a network—which is a long and tedious process—they typically use network measures and visualization techniques

to confirm or generate new hypotheses [80]. Visualization let them unfold the structure of their data, revealing potentially interesting social patterns between actors of the network. Analytics and visualization systems for SNA typically allow exploration of such data with the help of interaction, network measures, and data mining capabilities such as clustering directly implemented in the interfaces. Yet, most HSNA studies only give qualitative description of their network—which Rollinger call “soft” or “informal” network research—probably due to usability and formalism issues [2]. The coupling of visualization and data mining through interaction to support the generation of knowledge has been described as VA and can therefore provide support to social historians for their network construction, but also to go beyond simple qualitative description of their data. In this chapter, I first present a general overview of the field of visualization in §2.2 to share its utility and potential for social history. Then, I present the social history discipline with its use of quantitative methods in §2.3, before describing in depth how network analysis has been applied in the field in §2.4. Finally, I present in §2.5 how visualization and VA have been used in the context of HSNA, along the most popular systems currently used by social scientists and their limitations.

2.1 Notations

Social networks model social relationships or interactions between actors based on graphs. I note a graph

$$G = (V, E) \quad (2.1)$$

with V a set of nodes and $E \in V^2$ a set of edges. A weight is sometimes associated to edges, to quantify the strength or frequency of ties. In that case, each edge $e = (u, v, w)$ has a weight w with $u, v \in V$, $e \in E$, and $w > 0$. Multivariate graphs allow to associate attributes (also called properties) on nodes and edges. In that case, nodes and edges have associate attributes, such that a node $u = (u, v, a)$

$$a_u = (a_i, \dots, a_n) \quad (2.2)$$

with a_i, \dots, a_n the attributes of the node u defined on their domains A_i, \dots, A_n . We do not impose constraints for person nodes, but document nodes always have a time and location such that when $b = \text{document}$ then

Constructing a network from historical documents, which can vary tremendously in their formats and structures, is not a trivial task [2].

The most straightforward and well-known approach consists in constructing a network based on a simple graph such as in Equation 2.6, where the nodes are the persons extracted from the documents and a link is created when a mention of a social relationship is mentioned in the document (or if they appear in the same document) [?, 80]. This allows to have a simple network to visualize and analyze, but it does not always reflect the sociological complexity of information contained in the documents. HSNA network models have evolved over time to better take into account concrete properties of social networks, such as the importance of actors

or relations with weighted networks, multiple relationships with multiplex networks, dynamics of relations with dynamic networks.

Weighted networks model the importance of relations, with a weight w attributed to each edge $e = (u, v, w)$, with $u, v \in V$, $e \in E$, and $w > 0$. Multiplex networks allow to model multiple kinds of relationships between actors, such as spouses and witnesses relations for an historical network constructed from marriage acts. In that case, each edge $e = (u, v, d)$ of the graph

$$G = (V, E, D) \quad (2.3)$$

have a type $d \in D$ which characterize the relation. In the example of marriages, $D = \{spouse, witness\}$. Most relations extracted from historical documents also often contain time information, which can be modeled into dynamic networks. Many dynamic network models have been proposed [?], depending if the time is encoded in the nodes, the links, or both, and if entities have a discrete or interval time. As it is often hard to infer the end of social relationships from the trace of historical documents, we only consider in this thesis models which give a timestamp to either nodes or edges, such that

$$G = (V, E, T) \quad (2.4)$$

with vertices consisting of tuples (u, t) and edges of triples (u, v, t) , with $t \in T$.

Bipartite networks have been proposed to model relations between two types of entities, such as organization and employees where the relations link employees to organizations but not employees to employees or organizations to organizations [?]. Formally, each node of the graph

$$G = (V, E, B) \quad (2.5)$$

have a type $b \in B$, with $\text{card}(B) = 2$. For each edge $e = (u, v) \in E$, the types b_u and b_v of u and v are not equal $b_u \neq b_v$. Many social situations or documents can be modeled in these terms (affiliation lists or co-authoring). Multivariate networks, i.e., graphs, where vertices and edges can be assigned multiple “properties” or “attributes”, are less used in SNA. These attributes are often considered secondary, the emphasis of SNA being on the topology, its features, measures, and evolution.

Historians, demographers, sociologists, and anthropologists have also been designing specific data models for their social networks, based on genealogy or more generally kinship [59]. For genealogy, the standard GEDCOM [47] format models a genealogical graph as a bipartite graph with two types of vertices: individuals and families. This format also integrates an “event” object but it is diversely adapted in genealogical tools. The Puck software [58] has extended its original genealogical graph with the concept of “relational nodes” to adapt the data model to more family structures and to integrate other social relationships for anthropology and historical studies.

When creating a network, sociologists and anthropologists can use direct observations of the real world, which is not the case for historians who only have access to biased and partial sources. Indeed, the documents historians inspect are often produced by the political and economical elite of the time, and include the subjective view of the authors, especially for

literary sources (letters, journals, books, etc.). Historians therefore need to take a critical view on the sources by acknowledging the position of the authors of the documents compared to the rest of the society, and include it in the analysis [81]. Furthermore, the partiality of the sources often do not allow to have access to all possible relationships types of individuals. For example, if many formal relations can be extracted from official documents such as marriage acts and census, informal relations such as friendships can exist without leaving any written trace [80]. Even for official relationships such as parents and witnesses, there are high chances for missing documents, which do not allow to make too general and finite claims, such as “X is always the case” or “XX is never the case” [?]. Social historians therefore have to take into account the partiality and ambiguity of their sources into their analysis, in order to avoid including the bias inherent to their data into their high level historical conclusions.

2.2 Visualization

Visualization is often defined as “the use of computer-supported, interactive, visual representations of data to amplify cognition” [19]. Graphically displaying data allows us to leverage our visual system to gain a better acquisition of knowledge, leading to better decision-making, communication, and potential discoveries. The field of visualization can be split in three sub domains: **Scientific visualization** focus on visualizing continuous physically based data such as weather, astrophysics, and anatomical data, sometimes produced with simulations whereas **Information Visualization** is centered around the visualization of discrete abstract data points, often multidimensional. **Visual Analytics** emerged later from Information Visualization by mixing data mining and more complex analysis process with traditional information visualization displays. I focus in this thesis on the two former branches of visualization, as social scientists use both information visualization and VA systems to gain insight on the structure of the networks they are studying.

2.2.1 Information Visualization

Information Visualization focus on displaying abstract data to amplify cognition and gain insight on real world phenomena [19]. History is filled with classical examples of visual data displays which helped understand better specific events, such as Minard’s map of Napoleon’s march in Russia [45], or Snow’s dot map of cholera cases in London which showed the proximity between street pumps and cholera infections [135]. If several examples of information visualization can be found thorough history, it mainly developed as a scientific field in the 1960s with Tukey’s work on data analysis and visualization [140] and Bertin’s publication of Semiology of graphics [8].

In this foundational work, Bertin described and organized the different visual elements usable in graphical information displays, and linked them to data features and relations types. An illustration of this work of categorization for network data is illustrated in Figure 2.1. Michael Friendly writes that “To some, this appeared to do for graphics what Mendeleev had done for the organization of the chemical elements” [46]. The development of computer science and the rise of hardware capabilities during the same time created a big need for data visualization. The amount of data stored increased exponentially [?] and descriptive statistics were not enough to

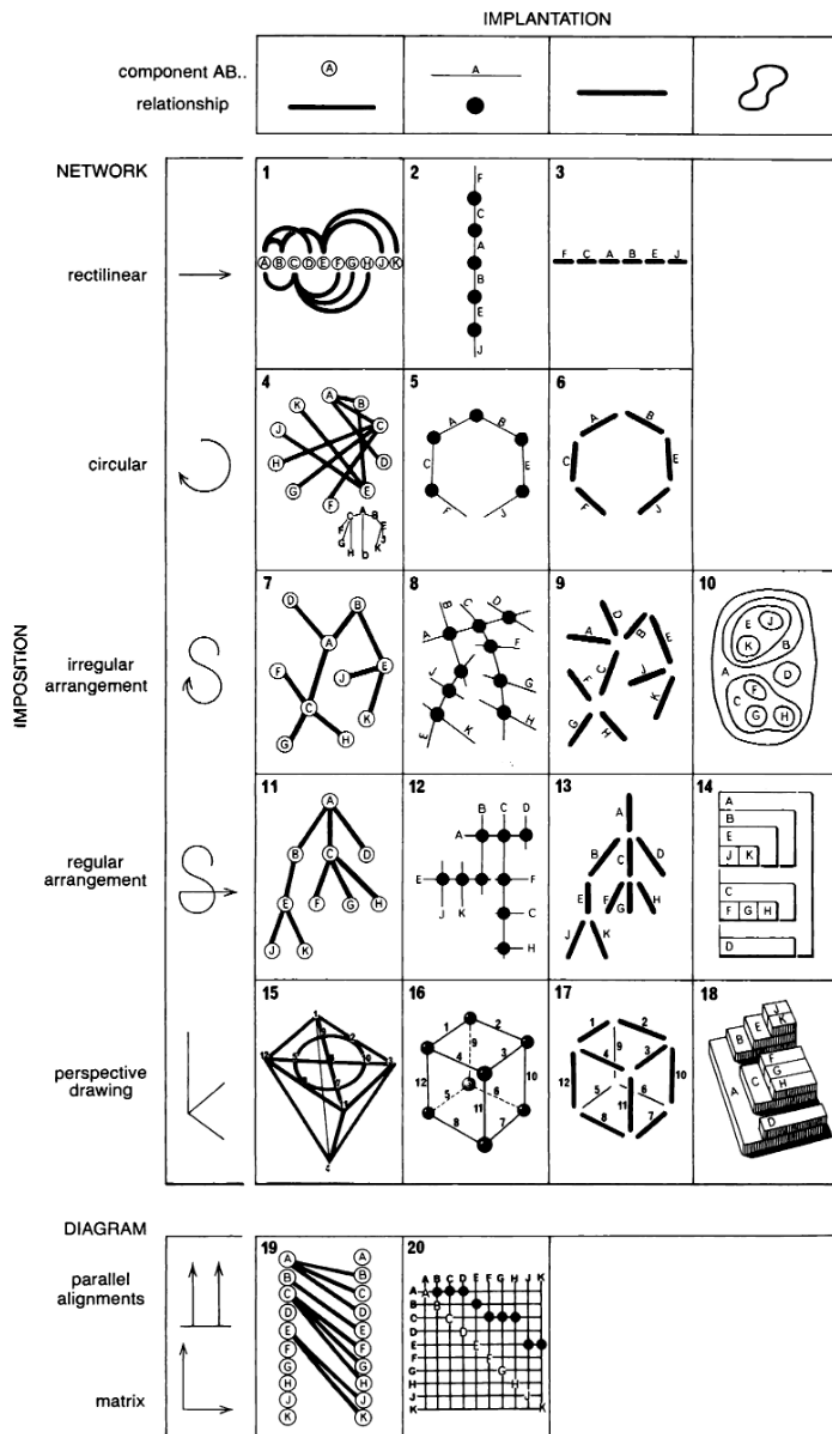


Figure 2.1 – Categorization of visual variables which can be used to represent network data, resulting in many different network representations. Image from [8].

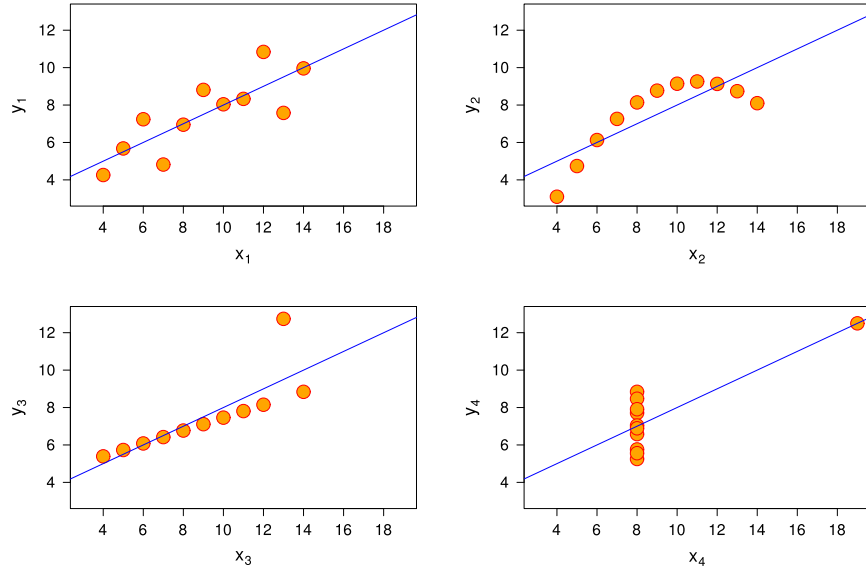


Figure 2.2 – Anscombe quartet. The four datasets have the same descriptive statistics (average, variance, correlation coefficient) but very different structures. Image from [4].

understand the underlying structure of the amount and diversity of produced data. Visualization, leveraging the human visual system, allowed to rapidly see the hidden structure of a dataset and detect interesting and unexpected patterns very often unseen with classical statistical methods. One classical illustration of this is Anscombe’s quartet [4] which consists of four datasets of 11 points in \mathbb{R}^2 with the same statistical measures (mean, variance, correlation coefficient, etc.) but with very different structures, that plotting the data show immediately. The four datasets are illustrated in Figure 2.2.

A large number of visualization techniques emerged to make sense of the diversity of data produced, such as multidimensional, temporal, spatial, or network data [?]. Instead of using taxonomies classifying graphics into categories such as histograms, pie charts, and stream graphs, some theorized how to describe graphics in a more systematic and structural way. In 1993, Wilkinson extended Bertin’s work and developed the Grammar of Graphics [?] as a way to describe the deep structure unifying every possible graphics, thus allowing to characterize and create graphics using common terms and rules. In this framework, a graphic can be defined as a function of six components: data (a set of data points and attributes from a dataset), transformations (statistical operations which modify the original data, e.g., mean and rank transformations), scales (e.g., linear and log scales), coordinate systems (e.g., cartesian and polar coordinate systems), elements (graphical marks such as rectangular or circular marks, and their aesthetics, e.g., color and size), and guides (additional information such as axes and legend). Many well-known visualization toolkits are now based on this framework, such as `vega` [125] and `ggplot` [?], as it allows greater expressiveness and reusability for graphic creation.

Visualization allows to gain insight on the structure of a given data, and has traditionally been used for confirmation and communication purposes, for example to verify hypothesis on empirical sciences, and later on to communicate findings, first to scientific peers, and nowadays to broader audiences for example through the means of data journalism [?].

2.2.2 Visual Analytics

VA consist in the coupling of visualization and data mining techniques to better support user in their knowledge generation process through continuous interaction with the data and statistical models [?]. It draws inspiration from exploratory visualization, interaction, and data mining. The process of exploratory visualization to gain new insights on the general structure of the data and potentially generate new hypotheses has been characterized by Tukey in 1960 as *Exploratory Data Analysis* [141]. It consists in trying to characterize the structure of a dataset with the help of continuous visualization and statistical measurements of different dimensions of the data. Visual exploration is enhanced by direct manipulation interfaces through interaction and usually follows the information-seeking mantra formalized by Schneiderman: “Overview first, zoom and filter, then details-on-demand” [?]. It allows users to first have a visual overview of the data and get an idea of its overall structure, to then change the point of focus to highlight interesting patterns with the help of filtering, querying, sorting, and zooming mechanisms. As the average size of datasets keeps growing, exploratory tools are often needed to make sense of large datasets and generate pertinent hypotheses.

More recent visual exploration interfaces also incorporate automatic analytical tools along with graphical displays, letting users apply data mining algorithms directly in the exploratory loop. This coupling of visualization and analytical methods such as data mining has been defined as Visual Analytics (VA) and is still a very active research field. Keim and al. define it as “a combination of automatic and visual analysis methods with a tight coupling through human interaction in order to gain knowledge from data” [71].

VA consist in the generation of knowledge using visualizations and statistical models of the data, that the user can explore using interaction. Such systems have been developed in various empirical domains, such as biology, astronomy, engineering, and social sciences, to explore various data types: multidimensional, temporal, geolocated, or relational (i.e., modeled into a network). Figure 2.3 shows the TULIP system, an example of a VA system developed for the analysis of network data. We discuss the uses of VA specifically for HSNA in §2.5.2

2.3 Quantitative Social History

Social History is a branch of history which aims at studying socio-economic aspect of past societies, with a focus on groups instead of specific individuals only. Charles Tilly writes that its goal is to “(i) documenting large structural changes, (2) reconstructing the experiences of ordinary people in the course of those changes, and (3) connecting the two” [?]. If the purpose of social history remained the same across time, methods and formalisms have evolved since its beginning in the 1930s. Specifically, the rise of computer science led to the development of quantitative history methods in the 1960s—now often referred as Digital Humanities—which brought new ways of grounding results in formalisms and quantitative models, instead of

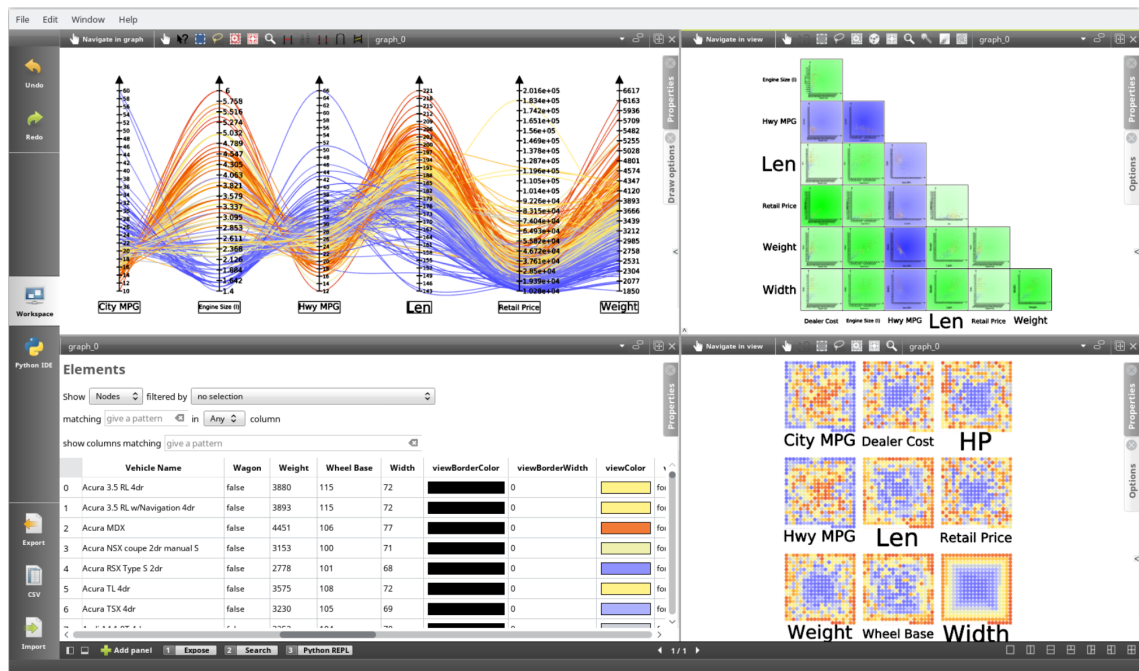


Figure 2.3 – TULIP software designed for application-independent network visual analytics [?]. The view shows a dataset among multiple interactive coordinated views. Users can also apply data mining algorithms on the data to extract interesting patterns.

solely relying on qualitative inspection of historical documents [?]. We discuss in this section the evolution of social history from the context of its beginning to the use of more recent quantitative approaches.

2.3.1 History, Social History, and Methodology

The concept of History is hard to define as its practice and codes highly evolved through time. Prost writes that "History is what historians do. The discipline called history is not an eternal essence, a Platonic idea. It is a reality that is itself historical, i.e. situated in time and space, carried out by men who call themselves historians and are recognized as such, received as history by various publics [113]." Retrospectively, History of a given time can thus be characterized by the different historical work produced at that time. Nevertheless, history can be characterized as the collection and study of historical documents to study and describe the past. As Langlois and Seignobos write, "The search for and the collection of documents is thus a part, logically the first and most important part, of the historian's craft" [?]. History emerged as a field with its own rules, conventions, and journals in the 1880s from faculties of letters, to counterbalance previous history works which were judged as too "literary" [?]. At that time and until now, two facets characterize the field, which are sometimes overlapping: one is political whereas the other one is methodological. The former aspect of history serves to create a shared story for the studied country and a sense of unity to its citizens. Antoine Prost says

that “it’s through history than France thinks itself” [113]. The latter aspect of history constitute a methodology to describe the past through methodical inspection of historical sources, in the aim of inferring dated facts about the past and trying to minimizing possible bias. Historical documents are thus at the core of the work of historians and having to cite historical documents and previous peers work to new claims is primordial to be considered as rigorous History work. However, methodological and epistemological facets (how historians should read and analyze their sources, how to cite them, what to report/not report, and what is the status a proof) of History have not been studied and discussed for a long time, until the end of the 1980s. Some historians were interested in historiography [18], but none were going to philosophical and epistemological reflexions of the History discipline. For Lucien Fèbvre, philosophising was even constituting a “capital crime” [40, 113].

Retrospectively, we can still observe shifts in the objects of study of historians through time, and their relation to sources. History was at first mainly event-centered and was focusing in characterizing central figures of the past like rulers and artists or shed light on central events like wars or political crisis. This narrative approach to history has been criticized for its open interpretation of historical documents, which can introduce bias from the authors [14]. In the 1930s, March Bloch and Lucien Fèbvre detached from traditional history by creating the “Annales school” (École des Annales) which aimed at placing the human as a component of a broader sociological, political, and economic system with influences between each other [16]. They strongly advised to exhaustively search from archives, to ground historical results in documents, texts and numbers. This new way of studying past events and societies became successful in a profession in crisis, by bringing a new lens of study on various societal subjects more grounded in sources and with a better intelligibility. This school of thought can be seen as one of the biggest milestones for Social History, which focuses on the socio-economical aspects of societies and their changes through time, rather than an event-centric view of History. For example, in his thesis, Ernest Labrousse—a well known figure of Social History—tries to describe and explain the economic crisis of France at the end of the “Ancien Régime”¹ through the evolution of the economic power of different social groups such as farmers, workers, property owners etc instead of solely describing memorable facts about the period [76]. Social History continued to evolve since the 1930s, introducing new methods and concepts, but always with the goals to describe periods and historical facts through a sociological lens and with a strong focus on sources and traceability.

2.3.2 Quantitative History

With the development of statistical methods and Computer Science, quantitative approaches of History emerged in the 1960s with the goal of analyzing numeric data directly extracted from historical documents. Economists led this first wave of quantification by studying past events using economical concepts and data. This approach, called “new economic history” or “cliometrics” was popularized by Fogel’s study on the economic impact of the development of railroads in America [?] and Fogel and Engerman’s controversial work on the

¹The “Ancien Régime” is an historical period of France which starts from the beginning of the reign of the Bourbon house at 1589 until the Revolution in 1789.

economy of slavery [?]. In the latter study, they extracted numbers of a sample of 5000 bills of slave sales from New Orleans to support the controversial claims that slavery was economically viable and that slaves had a decent material life, which brought up heated debate among the scientific community and the broad audience [?]. These kinds of approaches rapidly started to be used in other related domains such as demography, social history, and political history, sometimes rebranded as “new social history” and “new political history” [82] As extracting the data from raw documents and uploading it in computers—which were shared among whole departments—was very time-consuming at that time, “new history” projects often relied on a high division of labor among researchers, assistants, and students who operated with punch card operators [?]. Many saw the future of social sciences in computer programming, as Le Roy-Ladurie who wrote in 1968 “The historian of tomorrow will be a programmer, or he will not exist” [81].

However, quantitative methods started to be criticized in the 1980s with a vague of disillusionment, for several reasons. Stone was the first to raise his voice in 1979, after participating himself in several of those ambitious projects: “It is just those projects that have been the most lavishly funded, the most ambitious in the assembly of vast quantities of data by armies of paid researchers, the most scientifically processed by the very latest in computer technology, the most mathematically sophisticated in presentation, which have so far turned out to be the most disappointing” [?]. First, many researchers of this first wave dispensed themselves of source criticism, leading to simplification, anachronisms—such as using modern analytical categories and indices like the GDP—, and taking the numeric data from historical documents as objective. These problems could be in part explained by the fact that the work process was highly divided, meaning that the people analyzing the data did not necessarily inspect and read the original historical documents in depth. Secondly, the popularity of these methods made practitioners forget about the many biases inherent to statistics, such as the sampling bias, or the fact that historical data is essentially uncomplete data. This resulted in the computation of long data series and aggregates which were sometimes nonsensical given the gaps in the sources [81]. Finally, many historians raised their voice against the study of long-term trends instead of focusing on specific events and individuals. They challenged aggregations procedures and its assumptions, trying to go back to a more complex history by pointing that phenomena have to be studied and understood through several scales [?]. Indeed, computing correlations and aggregates at a national level greatly simplify complex phenomena, and misses specific group and individual related behaviours. Still, if their adoption remains slow and sometimes criticized among historians, quantitative methods provide tools to store, explore, and analyze historical documents systematically if used appropriately (i.e. not trying to bias the analysis, and not losing the trace of the original sources), especially that those methods highly evolved since the 1960s.

2.3.3 Digital Humanities

Digital Humanities is sometimes described as the second wave of computational social sciences [81]. The term has gained popularity since the 2010s and refer to “*research and teaching taking place at the intersection of digital technologies and humanities. Digital Humanities aims to produce and use applications and models that make possible new kinds of teaching*”

and research, both in the humanities and in computer science (and its allied technologies). Digital Humanities also studies the impact of these techniques on cultural heritage, memory institutions, libraries, archives and digital culture.” [?]. If the first wave of computational social sciences focused a lot on statistical methods such as regression models, correlation testing, and descriptive measures (mean, median, and variance) to make conclusions, digital humanities focuses more on the use of digital tools for exploration, teaching, and communication of humanities datasets and concepts, leveraging design, infographics, and interactive systems [?]. In the context of historical research, the term Digital History have been coined as “*an approach to examining and representing the past that works with the new communication technologies of the computer, the Internet network, and software systems. On one level, digital history is an open arena of scholarly production and communication, encompassing the development of new course materials and scholarly data collections. On another, it is a methodological approach framed by the hypertextual power of these technologies to make, define, query, and annotate associations in the human record of the past. To do digital history, then, is to create a framework, an ontology, through the technology for people to experience, read, and follow an argument about a historical problem.*” [?] Research which label itself as Digital History pivot around the curation and digitization of historical archives, the identification of historical concepts through computational and exploration methods, but also their communication to the general audience through digital technologies. Many Digital History projects are thus multidisciplinary by essence and involve several teams of researchers, such as the Republic of Letters project which consisted in digitizing, storing, and exploring letters of scholars across the world, in a common hub and using shared visualization tools [?]. It resulted in the elaboration of curated datasets and visualizations concerning the correspondence of various scholars such as Voltaire, Benjamin Franklin (see Figure 2.4), or John Locke, accessible in the same place by researchers and the general audience. With modern technologies and infrastructures, it also becomes possible to study large historical databases—often labeled under the term “big data”—as with the Venice Time Machine project [69] which aims at digitizing and analyzing thousands of documents from the archives of Venice to understand the political, geographical, and sociological dynamics of the cities across generations and centuries. Yet, some historians raised concern about this type of project, fearing that it could rapidly bring the same type of issues that we saw during the first wave of quantification, especially for big projects involving many actors and high ambitious goals [81].

Many projects which claim themselves as Digital History also leverage new methods compared to the 1960s and 1970s, such as the use of network methods and concepts [?]. Examples are the Viral Texts [?] and Living with Machines [?] projects which respectively study nineteenth-century newspapers and the industrial revolution by translating their sources into analyzable networks. We discuss more in depth the related work of network analysis for historical research in §2.4.

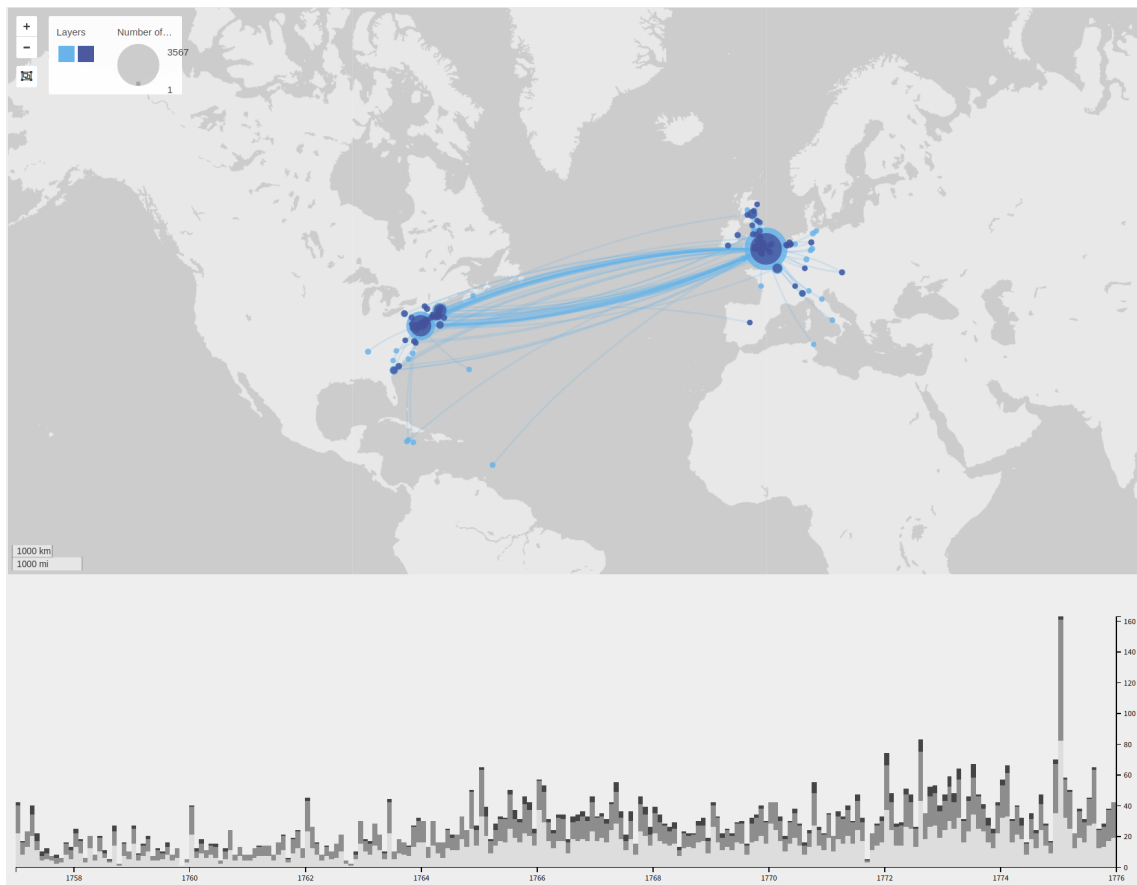


Figure 2.4 – Correspondence letters of Benjamin Franklin and his close relationships, using a map and an histogram, accessible online on the republic of letter website [?].

2.4 Historical Social Network Analysis

Historians started to use network analysis to study relational structures and phenomena of past societies in the 1980s, using similar methods developed by sociologist under the label of SNA. SNA is defined as an “*approach grounded in the intuitive notion that the patterning of social ties in which actors are embedded has important consequences for those actors. Network analysts, then, seek to uncover various kinds of patterns. And they try to determine the conditions under which those patterns arise and to discover their consequences*” [43]. the use of networks emerged in response to traditional sociology methods using pre-defined taxonomies and social categories to understand and explain sociological behaviors and phenomena, which could introduce bias. By modeling real observed social relationships and interactions with networks and by using mathematical and statistical methods to study those, sociologists have been able to explain sociological phenomena and describe sociological interactions through their direct observation modeled as networks. SNA is now a well-praised methodology in sociology and has been extended to historical research to study relational concepts such as kinship, friendships,

and institutions of the past. Social historians leverage their documents to extract relationships between entities—often persons—that they model into networks. Leveraging network measures and visualization, they can make conclusions through structural observations of such networks.

2.4.1 Sociometry to SNA

One of Sociology’s main goals is to study social relationships between individuals and find recurrent patterns and structures allowing to generalize on how social relations operate, and what are the social specificity of specific groups and individuals [127]. Traditional methods try to answer those questions using classical social classifications such as age, social status, profession, and gender, typically collected from surveys and interviews. Criticism pointed that this type of division is often partially biased and comes from predefined categories which are not always grounded in reality [43], and that using random sampling of individual with such methods remove them from their context. The sociologist Allen Barton wrote in this regard “For the last thirty years, empirical social research has been dominated by the sample survey. But as usually practiced, using random sampling of individuals, the survey is a sociological meatgrinder, tearing the individual from his social context and guaranteeing that nobody in the study interacts with anyone else in it” [?]. Sociometry is considered one of the bases of SNA and had the goal of redefining social categories through the lens of real social interactions and ties between persons, that sociologists wanted to observe in real conditions. It is in the 1930s that Moreno started to develop this new method by trying to depict real social interactions as a way to understand how groups and organizations were socially structured [95]. He developed sociograms to visually show friendships between people with the help of circles representing persons and lines modeling friendships. Figure 2.5 shows one of Moreno’s original sociograms to depict friendships in a class of first grades (left). Sociometry tremendously helped disseminate the metaphor of networks to model and understand social structures and phenomena. It was during the 1960s that sociologists and anthropologists took these concepts further and formalized SNA using graphs² and mathematical methods [?, 43], following the emergence of Graph Theory studies in the 1950s by Mathematicians such as Erdős [37]. Sociologists already had structural theories of social phenomena, and they rapidly saw the potential of graphs to model social relationships between actors of interest. Typically, a graph is noted

$$G = (V, E) \tag{2.6}$$

with V a set of vertices representing the actors of interest—typically persons—and $E \subseteq V^2$ a set of edges modeling social relationships. This simple model which do not take into account the diversity and extent of social relationships still allows the characterization of the sociological structure of groups and institutions—which is the primarily focus of SNA [43, 127]. More complex network models have been proposed with time to better take into account concrete properties of social relationships. I discuss those more in depth in §2.4.4.

²Graphs and networks refer to the same thing but are often used in different contexts. The term graph is preferred in a mathematical and abstraction setting, while the term network is mostly used when modeling real-world phenomena. We talk about nodes and links for networks and vertices and edges for graphs.

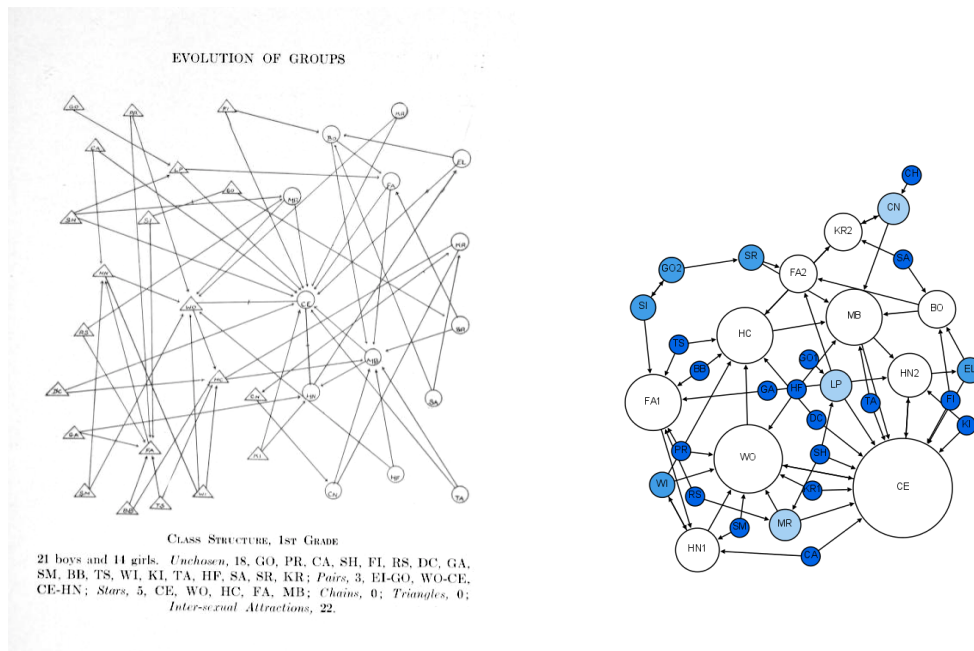


Figure 2.5 – Moreno’s original sociogram of a class of first grades from [94] (left). The diagram shows 21 boys (triangles) and 14 girls (circles). The same sociogram using modern practices generated from Gephi from [53] (right). The color encodes the number of incoming connections.

Graph theory brought a panoply of concepts and methods to characterize the structure of networks, that sociologists such as Coleman started to codify to use in a sociology setting [21]. The use of network measures let sociologists explain social phenomena through the formal description of real observations of relationships modeled as network.

2.4.2 Methods and Measures

Many measures and algorithms have been proposed in the network science and SNA literatures to characterize the structure of simple networks as defined in Equation 2.6 and relate it to social behaviours and phenomena [127, 138]. Network measures are either global or local, which allow to either make high level conclusions on the general structure of social relationships or individual behaviours. Widely used global measures include for example the density and the diameter, which give insight on the sparsity of the network and how distant on average are two random pairs of nodes. Conversely, local measures give information on the structural position of a node compared to the rest of the network. Centrality—probably the most used local measure—allows to formally compute a measure of how important or central are individuals in the network [?]. As defining what an important node is ambiguous, several types of centrality have been proposed such as the degree, betweenness, and closeness centrality, which respectively measure the number of connections, how nodes connects different groups, and how close are the nodes compared to the rest of the network.

More generally, sociologists aim at identifying recurring patterns of sociability between

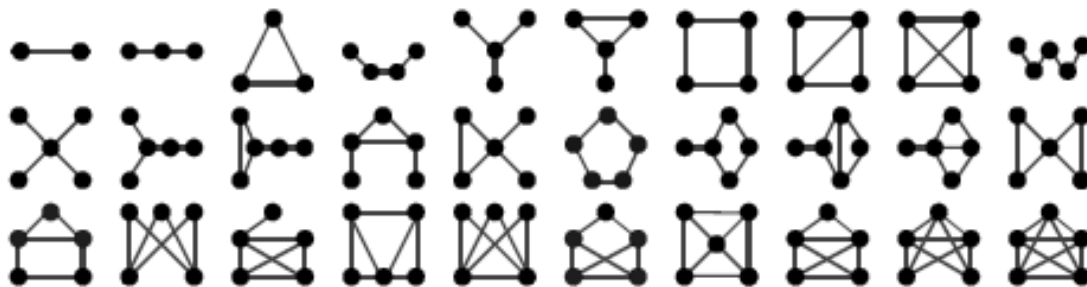


Figure 2.6 – All possible graphlets of size 2 to 5 for undirected graphs

actors, and link it to other behaviours, measures, or qualitative knowledge. These patterns can for example be small unconnected components, cliques, or bow-ties structures. Groups of nodes similarly located (central or distant) and having similar shapes are sometimes referred as structurally equivalent [80]. Instead of observing complex shapes, network scientists have also been interested at studying relationships at the lowest possible scale, i.e., observing relations between sets of 2 and 3 nodes at one, also called dyads and triads [148]. This reflects on Simmel's formal sociology, where he already referred to dyads and triads as the primal form of sociability [132]. More recently, graphlet analysis extended this concept to every pattern of N-entities [?].

Graphlets are defined as small connected *induced*, *non-isomorphic* subgraphs composing any network [91]. In an *induced* subgraph, two vertices linked in the original graph remain linked in the subgraph. For instance, if the original graph is a triangle \triangle we can only induce the simple edge $\bullet\text{---}\bullet$ or triangle \triangle subgraph (graphlet). The path of length 2 $\bullet\text{---}\bullet\text{---}\bullet$ has all vertices of the original graph but misses an edge and is, therefore, not a possible graphlet.

Figure 2.6 shows all graphlets of size 2 to 5, for undirected networks. Graphlets counting shows that graphlets are not found in a uniform distribution in social networks [?], thus revealing that social networks do not have the same structure that random networks. Precisely, entities in real-world networks tend to agglomerate into groups (also called clusters) where entities in the same groups interact more between them than with entities from other groups [?]. From a sociology perspective, it means that people tend to interact and socialize in groups and interact more rarely with other people from outside groups. These groups are often referred to as *communities*, and many algorithms have been proposed to find these automatically [?].

However, networks concepts, measures, and algorithms have not been used only to study groups, organizations, and societies, but also to focus on separate specific individuals. Indeed, two distinct methodologies emerged through the history of SNA: the structuralists and the school of Manchester [39, 43, 85].

The structuralists are interested in observing the relational structures and patterns forming a network, to make parallels between them and the social behaviors of actors in real life [79]. They think the positions of the persons in the network and the relational patterns they are part of reflect well the social activities and behavior in real life. Accordingly, sociologists in this school usually study organizations and specific groups—such as institutions, companies,

families—and want to explain their functioning through the description of the internal shapes and structures of the networks. Thus, they try to construct networks that exhaustively model all the interactions between the actors constituting the groups, as missing links would misrepresent the reality of interactions.

In contrast, the school of Manchester constituted by anthropologists focuses on studying specific individuals and all their interactions in the different facets of their lives and through time. They typically want to explain certain behaviors and social characteristics of individuals by their relationships and interactions in all their complexity and highlight the influence of different social aspects between them in one's life. One famous example is Mayer's study on austral Africa rural migrants going to cities [86] where he showed that the integration of urban mores and customs were directly correlated to the persons' relationships networks in the city. Xhosa³ people still interacting with rural people of their village in the city were less changing their customs. This school of thought typically relies on the concept of ego and multiplex networks [39]. Ego networks are networks modeling all the direct relations of one central node—in this case, a person—including the relations existing between the persons of this small network. They typically try to model the different types of relationships of a person, like their family, work, and friendship ties, and study them through time. By studying the ego network structure of someone, sociologists of this school try to leverage explanations on other social aspects of the persons like their social status, job, and gender. It is also common to compare several ego networks to make correlations between the social relationships of individuals and other interesting social categories [?].

These two methodologies of SNA are often not exclusives and current studies are typically inspired by those two traditions. This is especially true in history where even if historians may want to describe exhaustively a group or institution of the past, they are almost always interested in specific individuals they study in depth.

2.4.3 Historical Social Network Analysis

History started to use concepts and methods from SNA in the 1980s [149] in response to quantitative history, and to develop historical approaches—like *Microstoria* [49]—that focus on the study of individuals and small groups through the lens of their interactions and relationships directly extracted from historical documents. Beforehand, historians were already describing and studying relational structures such as families and organizations with qualitative methods and with classical taxonomies, without necessarily studying the relational aspect of these concepts. Network research allowed them to model those relational entities more thoroughly using networks, thus allowing them to make new observations that it was not possible to make without taking into account the relational structure of these entities [23]. Since then, HSNA—a term coined by C. Wetherell in 1998 [149]—has been applied by historians to study multiple types of relationships, like kinship [58], political mobilization [84], administrative and economic patronage [97]. If these approaches fall under some of the same critics as quantitative history [80] like leading to trivial conclusions, it still led to classical work and interesting

³Xhosa people are an ethnic group living in South Africa and talking the Xhosa language. and studied

discoveries, such as the study of the rise of the Medici family in Florence in the 15th century by Padgett & Ansell [105], or Alexander & Danowski study on Cicero's personal communications [?]. In this work, they modeled the communication of Cicero into a network using 280 letters written by him between 68 B.C. and 42 B.C. It allowed them to study the relationships between knights and senators—which is a subject of interest in Roman history—and concluded that knight-knight interactions were very rare compared to senator-senator and senator-knight interactions. Cicero communication network is illustrated in Figure 2.7.

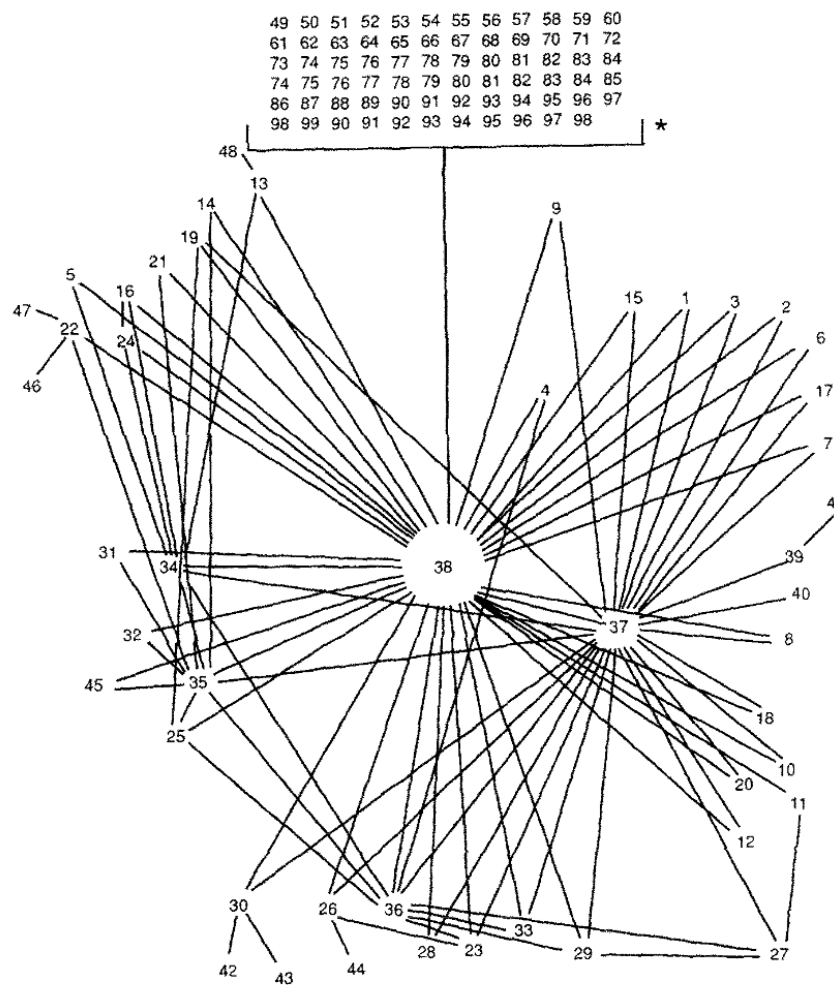


Figure 2.7 – Cicero personal communication network represented with a node-link diagram. Image from [?]

Several historians are using and continuously reflecting on HSNA methods [80] which can be very effective to study relational historical phenomena [72]. Moreover, historians rarely rely on a single approach when studying an era or phenomenon, they mix methods and tools from several domains of social and natural sciences with their own practices [105, 109].

2.4.4 Network Modeling

Constructing a network from historical documents, which can vary tremendously in their formats and structures, is not a trivial task [2].

The most straightforward and well-known approach consists in constructing a network based on a simple graph such as in Equation 2.6, where the nodes are the persons extracted from the documents and a link is created when a mention of a social relationship is mentioned in the document (or if they appear in the same document) [?, 80]. This allows to have a simple network to visualize and analyze, but it does not always reflect the sociological complexity of information contained in the documents. HSNA network models have evolved over time to better take into account concrete properties of social networks, such as the importance of actors or relations with weighted networks, multiple relationships with multiplex networks, dynamics of relations with dynamic networks.

Weighted networks model the importance of relations, with a weight w attributed to each edge $e = (u, v, w)$, with $u, v \in V$, $e \in E$, and $w > 0$. Multiplex networks allow to model multiple kinds of relationships between actors, such as spouses and witnesses relations for an historical network constructed from marriage acts. In that case, each edge $e = (u, v, d)$ of the graph

$$G = (V, E, D) \quad (2.7)$$

have a type $d \in D$ which characterize the relation. In the example of marriages, $D = \{spouse, witness\}$. Most relations extracted from historical documents also often contain time information, which can be modeled into dynamic networks. Many dynamic network models have been proposed [?], depending if the time is encoded in the nodes, the links, or both, and if entities have a discrete or interval time. As it is often hard to infer the end of social relationships from the trace of historical documents, we only consider in this thesis models which give a timestamp to either nodes or edges, such that

$$G = (V, E, T) \quad (2.8)$$

with vertices consisting of tuples (u, t) and edges of triples (u, v, t) , with $t \in T$.

Bipartite networks have been proposed to model relations between two types of entities, such as organization and employees where the relations link employees to organizations but not employees to employees or organizations to organizations [?]. Formally, each node of the graph

$$G = (V, E, B) \quad (2.9)$$

have a type $b \in B$, with $\text{card}(B) = 2$. For each edge $e = (u, v) \in E$, the types b_u and b_v of u and v are not equal $b_u \neq b_v$. Many social situations or documents can be modeled in these terms (affiliation lists or co-authoring). Multivariate networks, i.e., graphs, where vertices and edges can be assigned multiple “properties” or “attributes”, are less used in SNA. These attributes are often considered secondary, the emphasis of SNA being on the topology, its features, measures, and evolution.

Historians, demographers, sociologists, and anthropologists have also been designing specific data models for their social networks, based on genealogy or more generally kinship [59].

For genealogy, the standard GEDCOM [47] format models a genealogical graph as a bipartite graph with two types of vertices: individuals and families. This format also integrates an “event” object but it is diversely adapted in genealogical tools. The Puck software [58] has extended its original genealogical graph with the concept of “relational nodes” to adapt the data model to more family structures and to integrate other social relationships for anthropology and historical studies.

When creating a network, sociologists and anthropologists can use direct observations of the real world, which is not the case for historians who only have access to biased and partial sources. Indeed, the documents historians inspect are often produced by the political and economical elite of the time, and include the subjective view of the authors, especially for literary sources (letters, journals, books, etc.). Historians therefore need to take a critical view on the sources by acknowledging the position of the authors of the documents compared to the rest of the society, and include it in the analysis [81]. Furthermore, the partiality of the sources often do not allow to have access to all possible relationships types of individuals. For example, if many formal relations can be extracted from official documents such as marriage acts and census, informal relations such as friendships can exist without leaving any written trace [80]. Even for official relationships such as parents and witnesses, there are high chances for missing documents, which do not allow to make too general and finite claims, such as “X is always the case” or “XX is never the case” [?]. Social historians therefore have to take into account the partiality and ambiguity of their sources into their analysis, in order to avoid including the bias inherent to their data into their high level historical conclusions.

2.5 Social Network Visualization

Practitioners of SNA and HSNA have always visually depicted their network data for validation, exploration, and communication, mostly using node-link diagrams. With the use of more complex network models and the increase in average network size and density, new visualization techniques have been proposed to represent the diversity of studied networks. Moreover, more and more social scientists are following exploratory approaches using Visual Analytics (VA) tools, to describe more in-depth their data and generate new interesting hypotheses, using interaction and exploration capabilities.

2.5.1 Graph Drawing

Sociologists rapidly saw the potential of graphically showing relationships between individuals, to better comprehend the underlying social structure and communicate their findings [42]. Moreno elaborated sociograms to visually show friendships among schoolchildren with circles and lines to respectively show children and friendships ties [94]. This type of representation—commonly called node-link diagram—is the most widely used in social sciences, as it is rapidly understandable and effective for small to medium-sized networks which are predominant in social sciences. Finding an optimal placement for the nodes is however not that simple as several metrics can be optimized depending on the desired drawing, such as the number of edge crossings, the variance of edge length, orthogonality of edges, etc [24, 75]. Figure 2.8 shows some of these metrics, synthesized by Kosara and al. [75]. In Figure 2.5 we can see the difference in

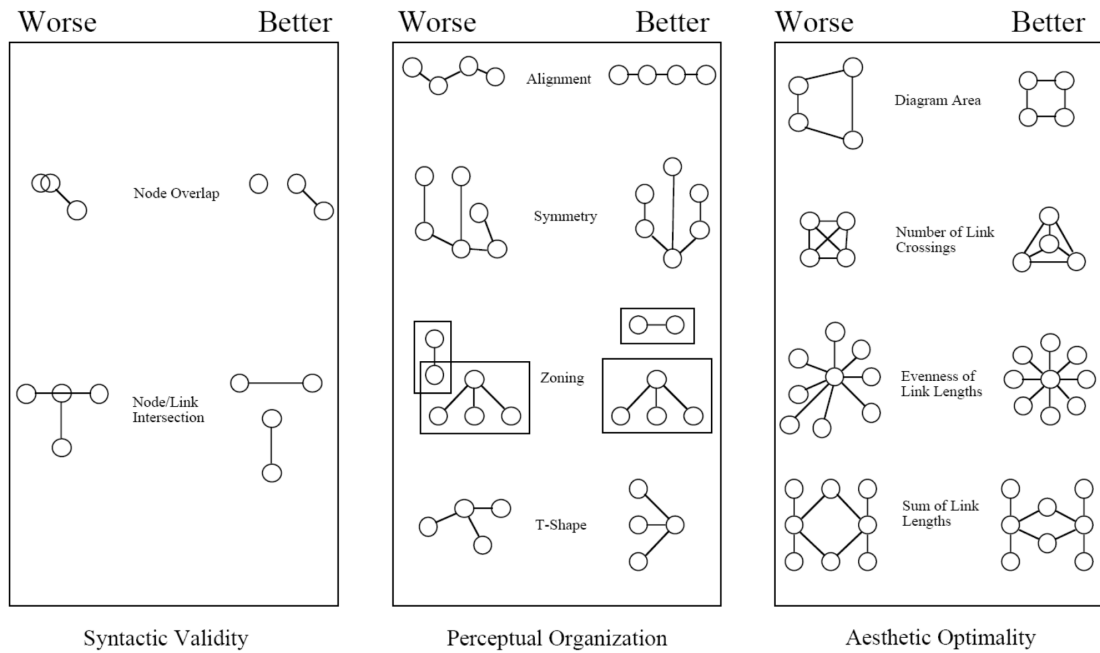


Figure 2.8 – Different criteria are proposed to enhance node-link diagram readability. Image from [75]

readability between the original manual layout (left) and an automatic one (right). Automatic layouts which aim at optimizing readability metrics give clearer diagrams. The number of edge crossings is often considered the most important measure, but finding a drawing with the optimal number of crossings is an NP-Hard problem, meaning that heuristics are needed for most real-world use cases. A large number of algorithms have been designed such as force-directed ones [?], modeling the nodes as particles that repulse each other and are attracted together when connected with a link that can be seen as strings. Other visual techniques have been proposed to represent networks such as matrices, circular layouts, and arcs, but are less used in social sciences [88]. Still, Matrices have been shown to be more effective than node-link diagrams for several tasks such as finding cluster-related patterns, especially for medium to large networks [?, 48].

As social scientists are using more complex network models such as bipartite or temporal networks, more sophisticated representations are needed. The visualization community developed new representations to visualize other network types such as dynamic hypergraphs with PAOHVis [142], clustered graphs with NodeTrix [64] (illustrated in Figure 2.9), geolocated social networks with the Vistorian [128], and multivariate networks with Juniper [102]. However, these new network representations take time to be adopted by social scientists who rarely use them.

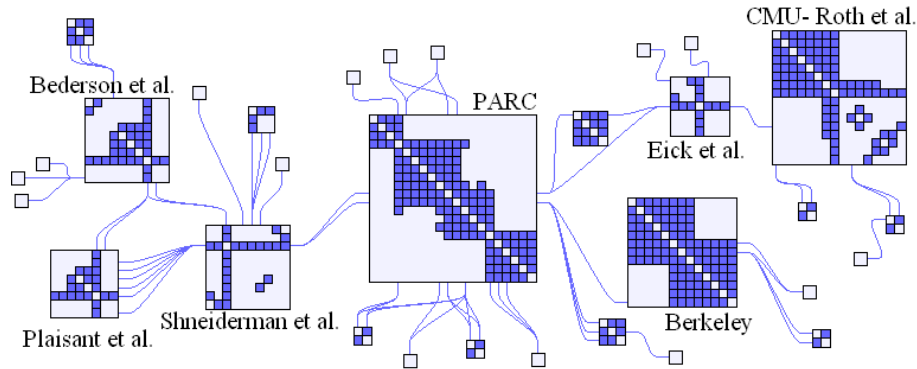


Figure 2.9 – NodeTrix system showing a scientific collaboration social network with clusters. Each cluster is represented as a matrix, Image from [64].

	Visualizations	SNA Measures and Models	Clustering	Filtering	Interaction/Direct Manipulation
Pajek	■■■	■■■	■■■	■■■	■■■
Ucinet	■■■	■■■	■■■	■■■	■■■
Gephi	■■■	■■■	■■■	■■■	■■■
NodeXL	■■■	■■■	■■■	■■■	■■■
Vistorian	■■■	■■■	■■■	■■■	■■■

Table 2.1 – Comparison table of most widely used visualization and analytical tool for HSNA. Visualizations: number of different visualization techniques, layout, and interactions. SNA and Models: Number of proposed SNA measures and algorithms. Clustering: Number of proposed clustering algorithms. Filtering: Possibilities of filtering according various criteria. Interaction/Direct Manipulation: Number of possible interactions mechanisms directly applicable on the visualizations.

2.5.2 Social Network Visual Analytics

Social scientists use visualization and analytical tools to gain insight on the structure of their finalized network data. Most widely used tools are Gephi [5], Pajek [98], Ucinet [?], and NodeXL [133] which provide node-link diagrams, implementations of network measures, algorithms, and clustering capabilities. Other SNA visualization tools have been proposed in the past such as Visone [?]. However, those softwares often do not include interaction and direct manipulation mechanisms, making the analysis more tedious for social historians and pose usability issues. In contrast, the Vistorian [128] let social historians visualize their network with multiple coordinated views (node-link, matrice, arc-diagram, and map), filters and direct manipulation, but do not integrate analytical options. Figure 2.10 shows the Vistorian interface used to explore an historical social network. I propose a classification of all those softwares in 2.1, which illustrate that most used software are analytical oriented (Pajek, Ucinet, Gephi, NodeXL) while the Vistorian is an interactive visualization tool. None of those software therefore fully correspond to the Visual Analytics label [?].

If analytical methods such as the computation of network measures, triad computation,

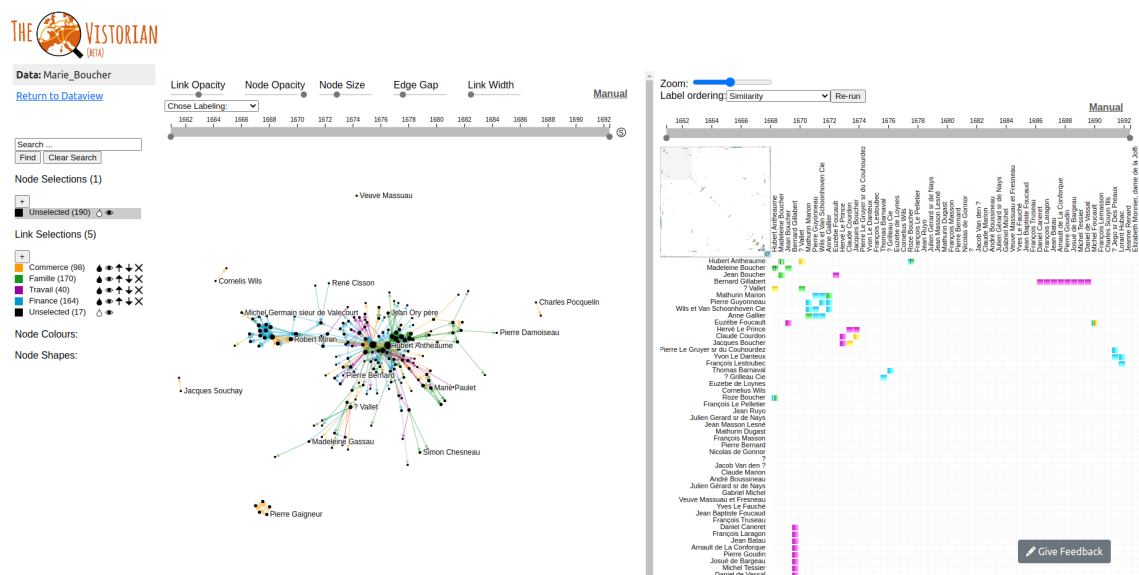


Figure 2.10 – Vistorian interface [128] used to explore a historical social network of business trades in the 17th century, with a coordinated node-link diagram and a matrice view.

or clustering provide a good framework to describe the structure of a network and link it to sociological explanation [127, 148], many social scientists such as historians are not trained in computer science and mathematical methods and therefore have trouble to first use those methods without guidance, but also interpret their results. This is particularly the case for black-box algorithms such as for clustering tasks: they usually end up trying several algorithms until they stumble upon a satisfactory enough solution [?].

Moreover, preparing and importing the data into visual and analytical software is complicated, as the annotation and network modeling process have not been globally formalized and every historian use different methods, formats and models. Many users do not succeed in importing their data in those systems without concrete help and guidance [2, 128] due to mismatches with data models, formats, or data inconsistencies (null values, white spaces, etc.). If they succeed visualizing their data, it often shows them these inconsistencies or errors such as duplications of entities or wrong attribute values. In other cases, they realize the network do not allow them to answer their sociological questions [80]. It leads to continuous back and forth between their analysis process inside the analysis tool they are using, and their annotation/modeling process, to correct errors or modify annotations. Interestingly, the network model choice plays a crucial role, as a simple network model representing only the persons (as it is often the case) makes it harder to trace back to the original documents containing the annotations from the network entities. Yet the majority of SNA systems enforce simple network models, making this retroactive process harder.

Some interfaces not primarily designed for social scientists incorporate data models encapsulating document representations, such as Jigsaw [136] which is a VA systems using textual

documents as a data model, originally developed for intelligence analysis. It allows an analysis of the documents and their mentions of entities (persons, locations, institutions, etc.) through multiple coordinated views. Using such model allow to rapidly see errors and inconsistencies in the documents annotations that the user can directly correct, while still following complex analyzes.

Finally, more work is still to be done on social network VA tools, to provide more guidance and power to social scientists while doing their analysis, and to help them to do easier back and forth between the annotation, network modeling, correcting, and analysis steps, as errors and inconsistencies can cause high variations in the network structure and hence the analysis results [31].

Bibliography

- [1] NodeXL: Simple network analysis for social media.
- [2] Mashael Alkadi, Vanessa Serrano, James Scott-Brown, Catherine Plaisant, Jean-Daniel Fekete, Uta Hinrichs, and Benjamin Bach. Understanding barriers to network exploration with visualization: A report from the trenches. *IEEE Transactions on Visualization and Computer Graphics*, 27(2), 2022. 13, 19, 21, 24, 40, 44, 45, 46, 48, 51, 53, 61
- [3] Keith Andrews, Martin Wohlfahrt, and Gerhard Wurzing. Visual Graph Comparison. In *2009 13th International Conference Information Visualisation*, pages 62–67, July 2009. 67
- [4] F. J. Anscombe. Graphs in Statistical Analysis. *The American Statistician*, 27(1):17–21, February 1973. 7, 28
- [5] Mathieu Bastian, Sebastien Heymann, and Mathieu Jacomy. Gephi: An open source software for exploring and manipulating networks. In Eytan Adar, Matthew Hurst, Tim Finin, Natalie S. Glance, Nicolas Nicolov, and Belle L. Tseng, editors, *ICWSM*. The AAAI Press, 2009. 18, 43, 65, 73
- [6] Sugato Basu, Ian Davidson, and Kiri Wagstaff. *Constrained Clustering: Advances in Algorithms, Theory, and Applications*. Chapman & Hall/CRC, first edition, 2008.
- [7] Leilani Battle and Jeffrey Heer. Characterizing Exploratory Visual Analysis: A Literature Review and Evaluation of Analytic Provenance in Tableau. *Computer Graphics Forum*, 38(3):145–159, 2019. 68
- [8] Jacques Bertin. *Sémiologie graphique: les diagrammes, les réseaux, les cartes*. Paris: Gauthier-Villars, 1967. 7, 26, 27
- [9] A. Bezerianos, F. Chevalier, P. Dragicevic, N. Elmqvist, and J.d. Fekete. GraphDice: A System for Exploring Multivariate Social Networks. *Computer Graphics Forum*, 29(3):863–872, 2010. 73
- [10] Marc Bloch. *Apologie Pour l'histoire*. A. Colin, 1949. 14
- [11] Christian Böhm and Claudia Plant. HISSCLU: A hierarchical density-based method for semi-supervised clustering. In *Proceedings of the 11th International Conference on Extending Database Technology: Advances in Database Technology*, EDBT '08, pages 440–451, New York, NY, USA, 2008. ACM.
- [12] Christian Bors, John Wenskovitch, Michelle Dowling, Simon Attfield, Leilani Battle, Alex Endert, Olga Kulyk, and Robert S. Laramée. A Provenance Task Abstraction Framework. *IEEE Computer Graphics and Applications*, 39(6):46–60, November 2019. 68

- [13] Michael Bostock, Vadim Ogjevetzky, and Jeffrey Heer. D³ Data-Driven Documents. *IEEE Transactions on Visualization and Computer Graphics*, 17(12):2301–2309, December 2011. 73, 83
- [14] Pierre Bourdieu. Sur les rapports entre la sociologie et l’histoire en Allemagne et en France. *Actes de la Recherche en Sciences Sociales*, 106(1):108–122, 1995. 31
- [15] Ulrik Brandes, Daniel Dellling, Marco Gaertler, Robert Gorke, Martin Hoefer, Zoran Nikoloski, and Dorothea Wagner. On Modularity Clustering. *IEEE Transactions on Knowledge and Data Engineering*, 20(2):172–188, February 2008. 52
- [16] Peter Burke. *History and Social Theory*. Polity, 2005. 31
- [17] Steven P. Callahan, Juliana Freire, Emanuele Santos, Carlos E. Scheidegger, Cláudio T. Silva, and Huy T. Vo. VisTrails: Visualization meets data management. In *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data - SIGMOD ’06*, page 745, Chicago, IL, USA, 2006. ACM Press. 53, 68
- [18] Charles-Olivier Carbonell. *L’Historiographie*. FeniXX, January 1981. 31
- [19] Stuart-K. Card, Jock-D. Mackinlay, and Ben Shneiderman. *Readings in Information Visualization: Using Vision to Think*. Morgan Kaufmann Publishers In, San Francisco, Calif, February 1999. 16, 26
- [20] Duen Horng Chau, Christos Faloutsos, Hanghang Tong, Jason I. Hong, Brian Gallagher, and Tina Eliassi-Rad. GRAPHITE: A Visual Query System for Large Graphs. In *2008 IEEE International Conference on Data Mining Workshops*, pages 963–966, December 2008. 67
- [21] J. S. Coleman. Introduction to mathematical sociology. *Introduction to mathematical sociology*, 1964. 36
- [22] TEI Consortium. TEI P5: Guidelines for electronic text encoding and interchange, February 2021. 51
- [23] Pascal Cristofoli. Aux sources des grands réseaux d’interactions. *Rezeaux*, 152(6):21–58, 2008. 13, 19, 38, 46, 48, 50, 52, 56, 57, 65
- [24] Pascal Cristofoli. Principes et usages des dessins de réseaux en SHS. *La visualisation des données en histoire*, page 35, 2015. 14, 41, 73
- [25] Pascal Cristofoli and Nicoletta Rolla. Temporalités à l’œuvre dans les chantiers du bâtiment. *Temporalités. Revue de sciences sociales et humaines*, (27), June 2018. 8, 49, 55, 68, 72
- [26] Tarik Crnovrsanin, Chris W. Muelder, Robert Faris, Diane Felmler, and Kwan-Liu Ma. Visualization techniques for categorical analysis of social networks with multiple edge sets. *Social Networks*, 37:56–64, 2014. 55

- [27] Gabor Csardi and Tamas Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*:1695, 2006. 67, 73
- [28] Erick Cuenca, Arnaud Sallaberry, Dino Ienco, and Pascal Poncelet. VERTIGo: A Visual Platform for Querying and Exploring Large Multilayer Networks. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2021. 67, 93
- [29] Zach Cutler, Kiran Gadhav, and Alexander Lex. Trrack: A Library for Provenance-Tracking in Web-Based Visualizations. In *2020 IEEE Visualization Conference (VIS)*, pages 116–120, October 2020. 80, 83
- [30] Mandeep K. Dhami, Ian K. Belton, and David R. Mandel. The “analysis of competing hypotheses” in intelligence analysis. *Applied Cognitive Psychology*, 33(6):1080–1090, 2019. 15
- [31] Jana Diesner, Craig Evans, and Jinseok Kim. Impact of Entity Disambiguation Errors on Social Network Properties. *Proceedings of the International AAAI Conference on Web and Social Media*, 9(1):81–90, 2015. 19, 45, 46, 50, 51
- [32] Dana Diminescu. The migration of ethnic germans from romania to west germany: Insights from the archives of the former communist regime. In *CERS, Public Lecture, UCLA*, Los Angeles, United States, March 2020. 49, 70
- [33] Nicole Dufournaud. La recherche empirique en histoire à l'ère numérique. *Gazette des archives*, 240(4):397–407, 2015. 13
- [34] Nicole Dufournaud. Comment rendre visible le rôle économique des femmes sous l'Ancien Régime ? Étude méthodologique sur les marchandes à Nantes aux XVIe et XVIIe siècles. In Bernard Michon and Nicole Dufournaud, editors, *Femmes et Négoce Dans Les Ports Européens (Fin Du Moyen Age - XIXe Siècle)*, pages 65–84. Peter Lang, 2018. 7, 16, 19, 48, 54
- [35] Nicole Dufournaud and Jean-Daniel Fekete. Comparaison d'outils pour la visualisation de sources historiques codées en XML/TEI. *Document numérique*, 9(2):37–56, April 2006. 51
- [36] Cody Dunne, Nathalie Henry Riche, Bongshin Lee, Ronald Metoyer, and George Robertson. GraphTrail: Analyzing large multivariate, heterogeneous networks while supporting exploration history. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '12*, pages 1663–1672, New York, NY, USA, May 2012. Association for Computing Machinery. 68
- [37] P. Erdős and A. Rényi. On the evolution of random graphs. In *On the Evolution of Random Graphs*, pages 38–82. Princeton University Press, October 2011. 35
- [38] Emily Erikson and Peter Bearman. Malfeasance and the Foundations for Global Trade: The Structure of English Trade in the East Indies, 1601–1833. *American Journal of Sociology*, 112(1):195–230, July 2006. 55

- [39] Michael Eve. Deux traditions d'analyse des reseaux sociaux. *Réseaux*, 115(5):183–212, 2002. 21, 37, 38
- [40] Lucien Febvre. VERS UNE AUTRE HISTOIRE. *Revue de Métaphysique et de Morale*, 54(3/4):225–247, 1949. 31
- [41] Jean-Daniel Fekete, Danyel Fisher, Arnab Nandi, and Michael Sedlmair. *Progressive Data Analysis and Visualization*. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, April 2019. 79
- [42] L. Freeman. Visualizing Social Networks. *J. Soc. Struct.*, 2000. 14, 16, 17, 41
- [43] L.C. Freeman. *The Development of Social Network Analysis: A Study in the Sociology of Science*. Empirical Press, 2004. 13, 15, 23, 34, 35, 37, 52
- [44] Manuel Freire, Catherine Plaisant, Ben Shneiderman, and Jen Golbeck. ManyNets: An interface for multiple network analysis and visualization. In *CHI '10*, CHI '10, pages 213–222, New York, NY, USA, 2010. ACM. 67
- [45] Michael Friendly. Visions and Re-Visions of Charles Joseph Minard. *Journal of Educational and Behavioral Statistics*, 27(1):31–51, March 2002. 26
- [46] Michael Friendly. A Brief History of Data Visualization. In Chun-houh Chen, Wolfgang Härdle, and Antony Unwin, editors, *Handbook of Data Visualization*, Springer Handbooks Comp.Statistics, pages 15–56. Springer, Berlin, Heidelberg, 2008. 26
- [47] GEDCOM: The genealogy data standard. 25, 41
- [48] Mohammad Ghoniem, J.-D. Fekete, and Philippe Castagliola. A comparison of the readability of graphs using node-link and matrix-based representations. In *IEEE Symposium on Information Visualization*, pages 17–24. Ieee, 2004. 42
- [49] Carlo Ginzburg and Carlo Poni. La micro-histoire. *Le Débat*, 17(10):133, 1981. 15, 38, 47
- [50] Barney G. Glaser and Anselm L. Strauss. *The Discovery of Grounded Theory: Strategies for Qualitative Research*. Aldine Transaction, New Brunswick, 5. paperback print edition, 2010. 15
- [51] Michael Gleicher. Considerations for visualizing comparison. *IEEE Trans. Vis. Comput. Graphics*, 24(1):413–423, 2018. 67
- [52] Claudia Goldin. Cliometrics and the Nobel. *Journal of Economic Perspectives*, 9(2):191–208, June 1995. 14
- [53] Martin Grandjean. Social network analysis and visualization: Moreno’s Sociograms revisited, 2015. 7, 36

- [54] Martin Grandjean. Analisi e visualizzazioni delle reti in storia. L'esempio della cooperazione intellettuale della Società delle Nazioni. *ME*, (2/2017), 2017. 64
- [55] Maurizio Gribaudo and Alain Blum. Des catégories aux liens individuels : l'analyse statistique de l'espace social. *Annales*, 45(6):1365–1402, 1990. 13
- [56] Jo Guldi and David Armitage. *The History Manifesto*. Cambridge University Press, October 2014. 13
- [57] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using NetworkX. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman, editors, *Proceedings of the 7th Python in Science Conference*, pages 11–15, Pasadena, CA USA, 2008. 67
- [58] Klaus Hamberger, Cyril Grange, Michael Houseman, and Christian Momon. Scanning for patterns of relationship: Analyzing kinship and marriage networks with Puck 2.0. *The History of the Family*, 19(4):564–596, October 2014. 25, 38, 41, 55, 57
- [59] Klaus Hamberger, Michael Houseman, and R. White, Douglas. Kinship network analysis. In John Scott & Peter J. Carrington, editor, *The Sage Handbook of Social Network Analysis*, pages 533–549. Sage Publications, 2011. 13, 25, 40
- [60] Mountaz Hascoët and Pierre Dragicevic. Interactive graph matching and visual comparison of graphs and clustered graphs. In Genny Tortora, Stefano Levialdi, and Maurizio Tucci, editors, *AVI '12*, pages 522–529. ACM, 2012. 67
- [61] J. Heer and D. Boyd. Vizster: Visualizing online social networks. In *IEEE Symposium on Information Visualization, 2005. INFOVIS 2005.*, pages 32–39, October 2005. 74
- [62] Jeffrey Heer. Agency plus automation: Designing artificial intelligence into interactive systems. *Proceedings of the National Academy of Sciences*, 116(6):1844–1850, 2019.
- [63] Louis Henry and Michel Fleury. Des registres paroissiaux a l'histoire de la population: Manuel de dépouillement et d'exploitation de l'état civil ancien. *Population (French Edition)*, 11(1):142–144, 1956. 14
- [64] Nathalie Henry, Jean-Daniel Fekete, and Michael J. McGuffin. NodeTrix: A Hybrid Visualization of Social Networks. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1302–1309, November 2007. 7, 42, 43
- [65] Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia D'amato, Gerard De Melo, Claudio Gutierrez, and Sabrina Kirrane et al. Knowledge graphs. *ACM Comput. Surv.*, 54(4), July 2021. 55
- [66] Pat Hudson and Mina Ishizu. *History by Numbers: An Introduction to Quantitative Approaches*. Bloomsbury Publishing, November 2016.
- [67] Infovis SC policies FAQ.

- [68] Piers J Ingram, Michael PH Stumpf, and Jaroslav Stark. Network motifs: Structure does not determine function. *BMC Genomics*, 7:108, May 2006. 66
- [69] Frédéric Kaplan. The Venice Time Machine. In *Proceedings of the 2015 ACM Symposium on Document Engineering*, DocEng '15, page 73, New York, NY, USA, September 2015. Association for Computing Machinery. 18, 33
- [70] Karine Karila-Cohen, Claire Lemerrier, Isabelle Rosé, and Claire Zalc. Nouvelles cuisines de l'histoire quantitative. *Annales. Histoire, Sciences Sociales*, 73(4):773–783, December 2018. 13, 47, 48, 53
- [71] Daniel Keim, Gennady Andrienko, Jean-Daniel Fekete, Carsten Görg, Jörn Kohlhammer, and Guy Melançon. Visual Analytics: Definition, Process, and Challenges. In Andreas Kerren, John T. Stasko, Jean-Daniel Fekete, and Chris North, editors, *Information Visualization: Human-Centered Issues and Perspectives*, Lecture Notes in Computer Science, pages 154–175. Springer, Berlin, Heidelberg, 2008. 7, 18, 29
- [72] Florian Kerschbaumer, Linda von Keyserlingk-Rehbein, Martin Stark, and Marten Düring. *The Power of Networks. Prospects of Historical Network Research*. Routledge, December 2021. 13, 14, 23, 39
- [73] Steffen Klamt, Utz-Uwe Haus, and Fabian Theis. Hypergraphs and cellular networks. *PLoS computational biology*, 5(5):e1000385, 2009. 92
- [74] Elena V. Konstantinova and Vladimir A. Skorobogatov. Application of hypergraph theory in chemistry. *Discrete Mathematics*, 235(1-3):365–383, May 2001. 92
- [75] C. Kosak, J. Marks, and S. Shieber. Automating the layout of network diagrams with specified visual organization. *IEEE Transactions on Systems, Man, and Cybernetics*, 24(3):440–454, March 1994. 7, 41, 42
- [76] Ernest Labrousse. *La Crise de l'économie Française à La Fin de l'Ancien Régime et Au Début de La Révolution*, volume 1. Presses Universitaires de France-PUF, 1990. 31
- [77] Charles-Victor Langlois and Charles Seignobos. *Introduction aux études historiques*. ENS Éditions, February 2014.
- [78] Katherine A. Larson. Thomas F. Tartaron, Maritime Networks in the Mycenaean World. New York: Cambridge University Press, 2013. *Comparative Studies in Society and History*, 56(4):1064–1065, October 2014. 13
- [79] Emmanuel Lazega. *Réseaux sociaux et structures relationnelles*. Presses universitaires de France, Paris, 1998. 17, 37
- [80] Claire Lemerrier. 12. Formal network methods in history: Why and how? In Georg Fertig, editor, *Social Networks, Political Institutions, and Rural Societies*, volume 11, pages 281–310. Brepols Publishers, Turnhout, January 2015. 13, 19, 24, 26, 37, 38, 39, 40, 41, 44, 45, 52, 64

- [81] Claire Lemerancier and Claire Zalc. *Quantitative Methods in the Humanities: An Introduction*. University of Virginia Press, March 2019. 14, 15, 19, 26, 32, 33, 41, 46, 47, 48, 53, 65
- [82] Claire Lemerancier and Claire Zalc. Back to the Sources: Practicing and Teaching Quantitative History in the 2020s. *Capitalism*, 2(2):473–508, 2021. 19, 32, 45, 46, 47, 48, 53
- [83] Bernard Lepetit. L'histoire quantitative : deux ou trois choses que je sais d'elle. *Histoire & Mesure*, 4(3):191–199, 1989. 13, 47
- [84] Carola Lipp. Kinship Networks, Local Government, and Elections in a Town in Southwest Germany, 1800-1850. *Journal of Family History*, 30(4):347–365, October 2005. 38
- [85] Gribaudo Maurizio. *Espaces, Temporalités, Stratifications : Exercices Méthodologiques Sur Les Réseaux Sociaux*. Editions de l'Ecole des Hautes Etudes en Sciences Sociales, Paris, January 2000. 37
- [86] Philip Mayer. Migrancy and the Study of Africans in Towns. *American Anthropologist*, 64(3):576–592, 1962. 38
- [87] Fintan McGee, Benjamin Renoust, Daniel Archambault, Mohammad Ghoniem, Andreas Kerren, and Bruno Pinaud et al. *Visual Analysis of Multilayer Networks*. Synthesis Lectures on Visualization. Morgan & Claypool Publishers, 2021. 55
- [88] Michael J. McGuffin. Simple algorithms for network visualization: A tutorial. *Tsinghua Science and Technology*, 17(4):383–398, August 2012. 42
- [89] Pierre Mercklé and Claire Zalc. Peut-on modéliser la persécution ? : Apports et limites des approches quantifiées sur le terrain de la Shoah. *Annales. Histoire, Sciences Sociales*, 73(4):923–957, December 2018.
- [90] R. Michalski, P. Kazienko, and D. Krol. Predicting Social Network Measures Using Machine Learning Approach. In *2012 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 1056–1059, Istanbul, August 2012. IEEE.
- [91] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. Network Motifs: Simple Building Blocks of Complex Networks. *Science*, 298(5594):824–827, October 2002. 37, 66
- [92] Christoph Molnar. *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. Lulu.com, 2019.
- [93] Megan Monroe, Krist Wongsuphasawat, Catherine Plaisant, Ben Shneiderman, Jeff Millstein, and Sigfried Gold. Exploring point and interval event patterns: Display methods and interactive visual query. *University of Maryland Technical Report*, 2012. 92

- [94] J. L. Moreno. *Who Shall Survive?: A New Approach to the Problem of Human Interrelations*. Who Shall Survive?: A New Approach to the Problem of Human Interrelations. Nervous and Mental Disease Publishing Co, Washington, DC, US, 1934. 7, 36, 41
- [95] J. L. Moreno. Foundations of Sociometry: An Introduction. *Sociometry*, 4(1):15, February 1941. 35
- [96] Zacarias Moutoukias. Buenos Aires, port between two oceans: Mobilities, networks, stratifications (2nd half of the 18th century). *E-SPANIA-REVUE ELECTRONIQUE D ETUDES HISPANIKES MEDIEVALES*, 25, 2016. 49, 70
- [97] Zacharias Moutoukias. Réseaux personnels et autorité coloniale : Les négociants de Buenos Aires au XVIIIe siècle. *Annales. Histoire, Sciences Sociales*, 47(4-5):889–915, October 1992. 38
- [98] Andrej Mrvar and Vladimir Batagelj. Analysis and visualization of large networks with program package Pajek. *Complex Adaptive Systems Modeling*, 4(1), April 2016. 43, 73
- [99] Natural earth. 73
- [100] Neo4j graph data platform. 65, 67, 83, 92
- [101] Rolla Nicoletta. Mobilité et conflits. Travailler sur les chantiers de construction piémontais dans la première moitié du XVIIIe siècle. In Andrea Caracausi and Marco Schnyder, editors, *Travail et Mobilité En Europe (XVIe-XIXe Siècles)*, Coll. Histoire et Civilisations. Presses universitaires du Septentrion, Villeneuve d'Ascq, 2018. 49
- [102] Carolina Nobre, Marc Streit, and Alexander Lex. Juniper: A Tree+Table Approach to Multivariate Graph Visualization. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):544–554, January 2019. 42
- [103] Gérard Noiriel. Naissance du métier d'historien. *Genèses. Sciences sociales et histoire*, 1(1):58–85, 1990.
- [104] Maryjane Osa. *Solidarity And Contention: Networks Of Polish Opposition*. Univ Of Minnesota Press, Minneapolis, first edition edition, July 2003. 13, 52
- [105] John F. Padgett and Christopher K. Ansell. Robust Action and the Rise of the Medici, 1400-1434. *American Journal of Sociology*, 98(6):1259–1319, May 1993. 7, 13, 17, 39
- [106] Pajek — Analysis and visualization of very large networks. 18, 20
- [107] Terence J. Parr and Russell W. Quong. ANTLR: A predicated-LL (k) parser generator. *Software: Practice and Experience*, 25(7):789–810, 1995. 83
- [108] Vanessa Peña-Araya, Tong Xue, Emmanuel Pietriga, Laurent Amsaleg, and Anastasia Bezerianos. HyperStorylines: Interactively untangling dynamic hypergraphs. *Information Visualization*, 21(1):38–62, January 2022. 59

- [109] Cindarella Sarah Maria Petz. *On Combining Network Research and Computational Methods on Historical Research Questions and Its Implications for the Digital Humanities*. PhD thesis, Technische Universität München, 2022. 14, 39, 47
- [110] James P. Philips and Nasseh Tabrizi. Historical Document Processing: Historical Document Processing: A Survey of Techniques, Tools, and Trends, September 2020.
- [111] Robert Pienta, Fred Hohman, Alex Endert, Acar Tamersoy, Kevin Roundy, Chris Gates, Shamkant Navathe, and Duen Horng Chau. VIGOR: Interactive Visual Exploration of Graph Query Results. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):215–225, January 2018. 67
- [112] Alexis Pister, Nicole Dufournaud, Pascal Cristofoli, Christophe Prieur, and Jean-Daniel Fekete. From Historical Documents To Social Network Visualization: Potential Pitfalls and Network Modeling. In *7th Workshop on Visualization for the Digital Humanities (VIS4DH)*, 2022. 45
- [113] Antoine Prost. *Douze Leçons sur l'histoire*. Média Diffusion, April 2014. 14, 23, 30, 31
- [114] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014. 67
- [115] Eric Ragan, Endert Alex, Jibonananda Sanyal, and Jian Chen. Characterizing Provenance in Visualization and Data Analysis: An Organizational Framework of Provenance Types and Purposes. *IEEE Transactions on Visualization and Computer Graphics*, 22(1), January 2016. 68
- [116] Ramana Rao and Stuart K. Card. The table lens: Merging graphical and symbolic representations in an interactive focus + context visualization for tabular information. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '94, pages 318–322, New York, NY, USA, 1994. Association for Computing Machinery.
- [117] Donghao Ren, Bongshin Lee, and Matthew Brehmer. Charticulator: Interactive Construction of Bespoke Chart Layouts. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):789–799, January 2019. 68
- [118] Pedro Ribeiro and Fernando Silva. Discovering Colored Network Motifs. In Pierluigi Contucci, Ronaldo Menezes, Andrea Omicini, and Julia Poncela-Casasnovas, editors, *Complex Networks V*, Studies in Computational Intelligence, pages 107–118, Cham, 2014. Springer International Publishing. 66
- [119] Giulio Rossetti and Rémy Cazabet. Community discovery in dynamic networks: A survey. *ACM Comput. Surv.*, 51(2), February 2018.
- [120] Fabrice Rossi, Nathalie Vialaneix, and Florent Hautefeuille. Exploration of a large database of French notarial acts with social network methods. *Digital Medievalist*, 9:2013, July 2014. 64

- [121] Juan A. Rubio-Mondejar and Josean Garrues-Irurzun. Women entrepreneurs and family networks in Andalusia (Spain) during the second industrial revolution. *Business History*, pages 1–22, May 2022. 13
- [122] C.J. Rueda and Catedral de Buenos Aires. *Matrimonios de La Catedral de Buenos Aires, 1747-1823*. Number v. 2 in Fuentes Históricas y Genealógicas Argentinas. Fuentes Históricas y Genealógicas Argentinas, 1989. 49
- [123] Anni Sairio. Methodological and practical aspects of historical network analysis: A case study of the Bluestocking letters. In Arja Nurmi, Minna Nevala, and Minna Palander-Collin, editors, *Pragmatics & Beyond New Series*, volume 183, pages 107–135. John Benjamins Publishing Company, Amsterdam, 2009. 55
- [124] Bahador Saket, Paolo Simonetto, and Stephen Kobourov. Group-level graph visualization taxonomy. In N. Elmqvist, M. Hlawitschka, and J. Kennedy, editors, *EuroVis - Short Papers*. The Eurographics Association, 2014.
- [125] Arvind Satyanarayan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. Vega-lite: A grammar of interactive graphics. *IEEE Trans. Vis. Comput. Graphics*, 23(1):341–350, 2016. 28, 83
- [126] Shrutika S. Sawant and Manoharan Prabukumar. A review on graph-based semi-supervised learning methods for hyperspectral image classification. *The Egyptian Journal of Remote Sensing and Space Science*, 2018.
- [127] John Scott. Social Network Analysis. *Sociology*, 22(1):109–127, February 1988. 23, 35, 36, 44, 52
- [128] Vanessa Serrano Molinero, Benjamin Bach, Catherine Plaisant, Nicole Dufournaud, and Jean-Daniel Fekete. Understanding the use of the vistorian: Complementing logs with context mini-questionnaires. In *Visualization for the Digital Humanities Workshop*, Phoenix, United States, October 2017. 7, 42, 43, 44, 51
- [129] Rachel Shadoan and Chris Weaver. Visual Analysis of Higher-Order Conjunctive Relationships in Multidimensional Data Using a Hypergraph Query System. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2070–2079, December 2013. 67
- [130] Termeh Shafie, David Schoch, Jimmy Mans, Corinne Hofman, and Ulrik Brandes. Hypergraph Representations: A Study of Carib Attacks on Colonial Forces, 1509-1700. *Journal of Historical Network Research*, pages 52–70 Pages, October 2017. 57, 64
- [131] Ben Shneiderman. Dynamic queries for visual information seeking. *IEEE Softw.*, 11(6):70–77, November 1994. 76
- [132] Georg Simmel. *Soziologie: Untersuchungen über die Formen der Vergesellschaftung*. Duncker & Humblot, Berlin, 7. aufl edition, 2013. 37

- [133] Marc A. Smith, Ben Shneiderman, Natasa Milic-Frayling, Eduarda Mendes Rodrigues, Vladimir Barash, Cody Dunne, Tony Capone, Adam Perer, and Eric Gleave. Analyzing (social media) networks with NodeXL. In John M. Carroll, editor, *Proceedings of the Fourth International Conference on Communities and Technologies, C&T 2009, University Park, PA, USA, June 25-27, 2009*, pages 255–264. ACM, 2009. 18, 43, 65, 73
- [134] SNA — Tools for social network analysis.
- [135] John Snow. On the Mode of Communication of Cholera. *Edinb Med J*, 1(7):668–670, January 1856. 26
- [136] John T. Stasko, Carsten Görg, and Zhicheng Liu. Jigsaw: Supporting investigative analysis through interactive visualization. *Inf. Vis.*, 7(2):118–132, 2008. 44, 58, 59
- [137] Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002.
- [138] Shazia Tabassum, Fabiola S. F. Pereira, Sofia Fernandes, and João Gama. Social network analysis: An overview. *WIREs Data Mining and Knowledge Discovery*, 8(5):e1256, 2018. 13, 17, 36
- [139] Natkamon Tovanich, Alexis Pister, Gaelle Richer, Paola Valdivia, Christophe Prieur, Jean-Daniel Fekete, and Petra Isenberg. VAST 2020 Contest Challenge: GraphMatchMaker: Visual Analytics for Graph Comparison and Matching. *IEEE Computer Graphics and Applications*, pages 1–1, 2021. 66, 67
- [140] John W. Tukey. The Future of Data Analysis. *The Annals of Mathematical Statistics*, 33(1):1–67, 1962. 26
- [141] John W. Tukey. *Exploratory Data Analysis*. Pearson, Reading, Mass, 1er édition edition, January 1977. 17, 29
- [142] Paola Valdivia, Paolo Buono, Catherine Plaisant, Nicole Dufournaud, and Jean-Daniel Fekete. Analyzing Dynamic Hypergraphs with Parallel Aggregated Ordered Hypergraph Visualization. *IEEE Trans. Visual. Comput. Graphics*, 27(1):1–13, January 2021. 42, 59, 93
- [143] Guido van Rossum. Python tutorial. Technical Report CS-R9526, Centrum voor Wiskunde en Informatica (CWI), Amsterdam, May 1995. 67
- [144] Ingeborg van Vugt. Using multi-layered networks to disclose books in the republic of letters. *Journal of Historical Network Research*, 1(1):25–51, October 2017. 55
- [145] Corinna Vehlow, Fabian Beck, and Daniel Weiskopf. The state of the art in visualizing group structures in graphs. In R. Borgo, F. Ganovelli, and I. Viola, editors, *Eurographics Conference on Visualization (EuroVis) - STARs*. The Eurographics Association, 2015.

- [146] VisMaster: Visual analytics — Mastering the information age.
- [147] Kiri Wagstaff, Claire Cardie, Seth Rogers, Stefan Schrödl, et al. Constrained k-means clustering with background knowledge. In *Icml*, volume 1, pages 577–584, 2001.
- [148] Stanley Wasserman and Katherine Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, November 1994. 17, 37, 44
- [149] Charles Wetherell. Historical Social Network Analysis. *Int Rev of Soc His*, 43(S6):125–144, December 1998. 13, 15, 23, 38, 46, 52, 54, 64
- [150] Kai Xu, Alvitta Ottley, Conny Walchshofer, Marc Streit, Remco Chang, and John Wenskovich. Survey on the Analysis of User Interactions and Visualization Provenance. *Computer Graphics Forum*, 39(3):757–783, June 2020. 53, 68
- [151] Michelle X. Zhou. “Big picture”: Mixed-initiative visual analytics of big data. In *Proceedings of the 6th International Symposium on Visual Information Communication and Interaction*, VINCI '13, page 120, New York, NY, USA, 2013. Association for Computing Machinery.