

1 **Interactive exploration of a global clinical network from a large breast cancer cohort**
2

3 Nadir Sella^{a,b,c,1}, Anne-Sophie Hamy^{b,d,e,2}, Vincent Cabeli^{c,3}, Lauren Darrigues^c, Marick Laé^{f,g}, Fabien
4 Reyal^{b,e}, Hervé Isambert^{c*}

5 a. Institut Roche, Boulogne-Billancourt, France

6 b. Residual Tumor & Response to Treatment Laboratory, RT2Lab, INSERM, U932 Immunity and
7 Cancer, Institut Curie, Paris, F-75248, France.

8 c. Laboratoire Physico Chimie Curie, Institut Curie, PSL Research University, CNRS UMR168,
9 75005, Paris, France

10 d. Department of Medical Oncology, Institut Curie, Saint-Cloud, F-92230, France.

11 e. Department of Surgery, Institut Curie, Université Paris, Paris, F-75248, France.

12 f. Department of Tumor Biology, Institut Curie, Paris, F-75248, France.

13 g. Department of Pathology, Henri Becquerel Cancer Center, INSERM U1245, UniRouen
14 Normandy University
15

16 **Corresponding Author information:** Hervé Isambert, herve.isambert@curie.fr, +33 1 56 24 64 74

17 **Co-corresponding Author information :** Pr Fabien REYAL, 0033144324087, 0033615271980,
18 fabien.reyal@curie.fr
19

20 **Abstract**

21 Despite unprecedented amount of information now available in medical records, health data remain
22 underexploited due to their heterogeneity and complexity. Simple charts and hypothesis-driven statistics
23 can no longer apprehend the content of information-rich clinical data. There is, therefore, a clear need
24 for powerful interactive visualization tools enabling medical practitioners to perceive the patterns and
25 insights gained by state-of-the-art machine learning algorithms. Here, we report an interactive graphical
26 interface for use as the front end of a machine learning causal inference server (MIIC), to facilitate the
27 visualization and comprehension by clinicians of relationships between clinically relevant variables. The
28 widespread use of such tools, facilitating the interactive exploration of datasets, is crucial both for data
29 visualization and for the generation of research hypotheses. We demonstrate the utility of the MIIC

30 interactive interface, by exploring the clinical network of a large cohort of breast cancer patients treated
31 with neoadjuvant chemotherapy (NAC). This example highlights, in particular, the direct and indirect
32 links between post-NAC clinical responses and patient survival. The MIIC interactive graphical
33 interface has the potential to help clinicians identify actionable nodes and edges in clinical networks,
34 thereby ultimately improving the patient care pathway.

35

36 **Keywords:** Machine learning - data visualization- residual cancer burden - neoadjuvant chemotherapy
37 – breast cancer.

38

39 INTRODUCTION

40

41 The availability of health data from patient medical records is increasing, and these data
42 constitute, in theory, a rich resource for research purposes. However, despite the unprecedented amount
43 of information now available, health data remain underexploited due to their heterogeneity and
44 complexity. There is, therefore, an urgent need for innovative tools, based on intuitive and interactive
45 graphical interfaces, specifically designed for the exploration of health data by medical practitioners.
46 Data visualization is gradually emerging as a new field of research, and graphical representations are
47 used for two main purposes: (i) explanatory illustration, to highlight novel scientific insights graphically
48 and to ensure efficient communication between scientists¹⁻⁴; and (ii) exploratory analysis, searching
49 for relationships previously overlooked and leading to new discoveries, thereby maximizing the
50 potential of information-rich databases. We present here an *exploratory analysis* of a global clinical
51 network from a large breast cancer cohort, with a novel interactive graphical interface for the exploration
52 of health data.

53 We previously developed an advanced computational method for graphical analyses, including
54 causal relationships, from multivariate data⁵. The underlying MIIC (Multivariate Information-based
55 Inductive Causation) algorithm, which was released as an online server⁶, uses a machine learning
56 method combining constraint-based and information theory approaches to reconstruct causal, non-causal
57 or mixed networks from large datasets. The MIIC algorithm was first developed to analyze categorical

58 genomic data^{5, 6} and has recently been extended to the analysis of more challenging heterogeneous
59 datasets, such as medical records, combining both categorical and continuous variables, in which
60 interdependence is notoriously difficult to assess⁷.

61 Breast cancer (BC) clinical datasets are particularly suitable for the type of exploratory analysis
62 presented here, as BC is a complex heterogeneous disease highly variable in its aggressiveness and
63 prognosis. BC remains one of the leading causes of cancer-related death among women. The BC
64 patients included in the cohort analyzed here were treated with neoadjuvant (or preoperative)
65 chemotherapy (NAC). NAC was originally restricted to patients with inflammatory or locally advanced
66 BC, but is now the standard care for aggressive early-stage breast cancers, *i.e.* triple-negative (TNBC)
67 and *HER2*-positive BCs^{8,9}. From the patient's viewpoint, the benefits of the neoadjuvant strategy include
68 a greater feasibility of breast-conserving surgery and the prognostic stratification of risk obtained after
69 analyses of the residual tumor burden at surgery. From the research and development standpoint, the
70 neoadjuvant setting makes it possible to monitor the chemosensitivity of the tumor *in vivo*, and provides
71 an opportunity for the rapid validation of research hypotheses and the acceleration of drug approval.

72

73 **RESULTS**

74 The global network displayed in Fig. 1 is accessible at
75 https://miic.curie.fr/job_results_showcase.php?id=NEOREP. We discuss below some of the links
76 inferred in the NEOREP network after grouping according to several clinically relevant concepts
77 identified from published studies on BC.

78

79 ***MIIC performs quality control***

80 MIIC first identifies relationships between a disease and the corresponding treatment. ER
81 positivity — which is predictive of efficacy for anti-hormonal treatment¹⁰ — is associated with the use
82 of endocrine therapy (Fig. S3A), and a similar association is observed for *HER2*-positivity and
83 trastuzumab use (Fig. S3B)¹¹. Beyond cancer, significant associations are also found between depression
84 and the use of psycholeptics (Fig. S3C), between thyroid disorders and thyroid hormone use (Fig. S3D),

85 and between hypertension and drugs for the treatment of cardiovascular diseases (Fig. S3E). More
86 generally, comedication use is associated with the type of NAC (Fig. S3F), reflecting the greater
87 likelihood of less toxic regimens being prescribed to fragile patients (patients on other types of
88 medication) than to patients without comedication¹²⁻¹⁴.

89 MIIC then identifies clinical factors known to be epidemiologically related (Fig. S4A). Menopause, a
90 process occurring in older women, is directly linked to age (Fig. S4B) (median age: 43 years for
91 premenopausal, versus 58 years for postmenopausal women). Postmenopausal status is associated with
92 dyslipidemia (Fig. S4C)¹⁵. Consistent with these associations, body mass index (BMI) increases with
93 age (Fig. S4A, S4D) and both factors, which have been reported to increase cardiovascular risks, are
94 linked to hypertension (Fig. S4A, S4E). The number of drugs taken by a patient (comedication) increases
95 with the number of comorbidities (Fig. S4A, S4F).

96

97 *MIIC identifies inherent associations between variables*

98 The duration of neoadjuvant treatment is directly linked to the type of NAC regimen delivered (Fig. 2A)
99 reflecting the fact that anthracycline-based (AC) regimens usually include four cycles (median of 106
100 days, Fig. 2B), whereas sequential regimens in which anthracyclines are followed by taxanes are
101 generally administered over six or eight cycles (median of 147 days, Fig. 2B). The number of nodes
102 retrieved is associated with the type of axillary surgery (Fig. 2C), consistent with the fact that sentinel
103 node (SLN) biopsy procedures were developed to reduce the number of lymph nodes removed during
104 dissection (LND) (Fig. 2D)¹⁶. MIIC correctly represents the direct links between residual cancer burden
105 (RCB) (Fig. 2E) and the patterns making up this score, derived from measurements on the primary tumor
106 bed (size, fraction of invasive cancer, cellularity) and the regional lymph nodes (number of positive
107 lymph nodes).

108

109

110 *MIIC identifies intra- and inter-modality associations*

111

112 For the variables derived from pathology records, MIIC found associations between tumor
113 grade, Ki67, and mitotic index (Fig. S5A-B-C), all of which are markers of tumor proliferation¹⁷. MIIC
114 can also visualize links between patterns assessed in different ways. Measurements of pre-NAC tumor
115 size evaluated clinically, by mammography and by MRI, were found to be closely related (Fig. S5C-E)
116 as previously reported^{18, 19}. Similarly, the response to treatment assessed clinically at NAC completion
117 was found to be associated with histological size based on the surgical specimen (Fig. S5F).

118

119 ***MIIC provides insight into tumor biology and response to treatment***

120 The presence of lymphovascular invasion (LVI) in the post-NAC specimen is associated with a
121 higher RCB index, consistent with the strong resistance to chemotherapy of these tumors²⁰ (Fig. S6A).
122 TNBCs and *HER2*-positive tumors have a higher pre-NAC mitotic index and more stromal TIL
123 infiltration (Fig. S6B-C) than luminal BCs^{21,22}. Consistently, high TIL levels are significantly associated
124 with histological grade 3 tumors (Fig. S6D).

125

126 ***MIIC reflects clinical practice***

127 Several associations highlighted in the network reflect clinical practice decisions applied
128 throughout BC centers. For example, the likelihood of performing conservative breast surgery depends
129 on tumor histology (higher rates of mastectomy have been reported for patients with lobular or other
130 histological types of tumor less likely to respond to NAC)^{23, 24} (Fig. S7A) and is positively associated
131 with the practice of oncoplastic surgery²⁵ (Fig. S7B). Similarly, lumpectomy is more frequently
132 associated with radiation therapy than with mastectomy (Fig. S7C)²⁶⁻²⁹. After surgery, the addition of a
133 second line of treatment by adjuvant chemotherapy, to decrease the risk of relapse, is driven by the
134 identification of factors associated with a poor prognosis³⁰, such as high levels of lymph node
135 involvement (Fig. S7D).

136 Beyond these well-established practices, MIIC also identified differences in clinical practices between
137 the two centers of the cohort (Fig. 3A). For example, oncoplastic surgery and adjuvant chemotherapy
138 were performed at only one of the two centers (Fig.3B-C); the NAC regimen also differed between

139 centers, with the Curie St Cloud center using more AC regimens than AC-taxane combinations, resulting
140 in a shorter duration of NAC treatment (Fig. 3D-E).

141

142 ***MIIC traces the natural course of the disease***

143

144 The natural course of BC may include local relapse, possibly followed by distant metastases,
145 the trigger events leading to death³¹⁻³⁵ (Fig. 4A-C). Contralateral BC is often used in composite survival
146 endpoints, such as distant relapse-free survival³⁶, but MIIC clearly identifies contralateral BC as an event
147 being independent of other oncologic events and almost totally isolated from the rest of the network
148 (Fig.1). Luminal BC is known to recur and develop metastases later than *HER2*-positive BC and TNBC
149 (Fig. 4D)^{21, 22, 37, 38}. The link between has also been found between PR negativity and a higher risk of
150 brain metastasis³⁹⁻⁴³ (Fig. 4E).

151

152

153 ***MIIC identifies unexpected associations, leading to new discoveries***

154 With more than 15 associations involving treatment center (Fig. 3A), MIIC unmasked an
155 unexpected “batch” effect relating to the site of BC treatment in this cohort. The observed differences
156 reflect not only differences in therapeutic practice, but also in the characteristics of the population
157 (differences in the proportion of women with psychological disorders, difference in incomes), in tumor
158 presentation (tumor size), in pathological variable scoring (grade, presence of pre-NAC LVI, tumor
159 cellularity, TILs), and in time to treatment within the care pathway.

160

161 ***MIIC identifies factors likely to improve prediction or prognosis***

162 MIIC also favors new insights, e.g. comedication appears to protect against local relapse (Fig.
163 5A). Several retrospective studies have reported this association, with the use of statins⁴⁴, NSAIDs⁴⁵, or
164 beta-blockers⁴⁶ found to have indirect anticarcinogenic effects. It has recently been suggested that these
165 non-oncological treatments may have immunomodulatory and chemosensitizing effects⁴⁷.

166

167 ***MIIC suggests relevant combinations of predictive of prognostic biomarkers***

168 MIIC may provide clues to combinations of new prognostic biomarkers likely to improve the
169 prediction of response to chemotherapy, or post-NAC prognosis. Pre-NAC lymphovascular invasion
170 (LVI) was found to be associated with both lower rates of clinical response (Fig. 5B) and shorter relapse-
171 free survival (Fig. 5C). Both RCB (Fig. 5D-E) and post-NAC mitotic index (Fig.5D-F), a parameter
172 rarely used in practice but nevertheless reported to be a predictor of BC recurrence^{48, 49}, appear to be
173 strongly associated with the risk of death. MIIC may, therefore, be an efficient tool for identifying
174 features likely to improve prognosis, by combining gold standard indicators with other parameters, such
175 as post-NAC mitotic index, and post-NAC LVI, for example. Finally, MIIC also makes it possible to
176 optimize the binning of residual cancer burden (RCB). RCB is a post-NAC histological score calculated
177 as an increasing continuous index, and then subdivided into four classes (0, I, II, and III)⁵⁰. Our analysis
178 based on information maximization principles suggested a new unsupervised classification of RCB
179 scores into three categories (Fig. 5E), with RCB=0 with low RCB values merged, in particular, into a
180 single class associated with a good prognosis.

181

182 **DISCUSSION**

183

184 When applied to a large cohort of BC patients, the MIIC algorithm successfully (i) performed
185 quality controls; (ii) identified intra- and inter-modality correlations; (iii) highlighted differences in
186 clinical practice, including center specificities; (iv) traced the natural course of the disease; (v)
187 highlighted unsuspected and hidden associations, leading to new discoveries. The interactive
188 visualization and causal analyses provided by this algorithm make it a promising tool for fast and
189 effective explorations of the increasing amount of available health data.

190 The amount of exploitable health data is increasing exponentially. The best known health data
191 resource for cancer studies remains the SEER (Surveillance, Epidemiology, and End Results) database,
192 which collects data from population-based cancer registries covering approximately 34.6% of the US
193 population^{51, 52}. By 2016, the National Cancer Database (NCDB) had amassed more than 34 million
194 hospital records from cancer patients (almost four times the size of the SEER database), to become the

195 largest clinical cancer registry in the world⁵³. In France, the French administrative health care database,
196 the SNDS (*Système National des Données de Santé*), is one of the largest administrative databases in
197 the domain of medicine, providing many opportunities for medical research^{54,55}, as it covers 99% of the
198 French population (about 66 million people). The French government is planning to ease access to this
199 almost exhaustive population research resource, through release as part of the “Health data hub” project.
200 Finally, beyond these structured databases, the largest mine of untapped data worldwide remains the
201 content of electronic health records (EHRs), encompassing a full range of data (clinical notes, laboratory
202 results, imaging, genetic data, etc.) relating to patient care. Recent advances in information technology
203 have made it easier for both hospitals and healthcare institutions to collect large amounts of healthcare
204 data.

205 Biomedical scientists are now facing new challenges in the management and analysis of
206 massive, heterogeneous datasets⁵⁶. These challenges include the development of tools for exploration
207 and visualization, analytical methods, integration into a comprehensive overview, and translation of the
208 findings into public health impact. The visualization of information makes it possible for users to find
209 profound patterns in clinical data, through visual recognition. Simple charts cannot represent the
210 complexity of big data analyses and fail to support multifaceted tasks effectively^{3,4}. There is, therefore,
211 a need for sophisticated visualization tools dealing with many elements simultaneously and enabling
212 users to perceive the patterns and insight generated by the algorithm⁵⁷. Supplementary Table 1 shows
213 the main data visualization tools used to present medical data. Many of the visual methods have been
214 adopted directly from the field of data mining, but others, specific to the healthcare domain, have also
215 been designed (Supplementary Table 2). For example, Happe and Drezen built the ePEPs toolbox, which
216 displays relevant patterns extracted by eye from patient reimbursement data in the SNDS database, and
217 supporting interactive exploration by researchers⁵⁸. CARRE provides web-based components for
218 interactive health data (fitness and biomarkers) visualization and risk analysis for the management of
219 cardiorenal diseases⁵⁹. The MITRE Corporation has also developed a web-based solution that provides
220 an overview of an individual’s health through graphical representations of EHR data, highlighting
221 abnormal values⁶⁰. None of these visualization programs has yet managed to bridge the gap between of

222 the large amounts of clinical data available and the discovery of clinical knowledge or paths for scientific
223 research. By processing large heterogeneous sets of variables inherent to clinical records, MIIC provides
224 physicians with a full picture of BC disease.

225 In addition to this use for visualization, the MIIC algorithm presents several other advantages
226 for analyses, including its unsupervised nature, overcoming the need for training or human involvement.
227 This feature makes it possible to obtain new knowledge through the automatic identification of patterns
228 and dependences in the data, highlighting new interactions, and it may be of use for feature selection in
229 machine learning models.

230 In conclusion, MIIC, an open-access, interactive, multitask tool, is designed to visualize datasets to help
231 clinicians and researchers to understand the relationships between the variables within them. It opens
232 up promising perspectives for guiding the generation of new hypotheses, helping clinicians identify
233 actionable nodes and edges in clinical networks, and revealing new clues to relationships of interest for
234 research purposes. Its widespread use in the field of health data could increase the accuracy of prediction
235 for treatment responses and prognosis. This tool has the potential to improve the care pathway and,
236 ultimately, the survival of patients.

237

238 **METHODS**

239 *Patients and treatment*

240 We analyzed a cohort of 1197 patients with non-metastatic BC treated by NAC, with or without
241 trastuzumab, followed by surgery, at either of the two Institut Curie sites (Paris and Saint Cloud)
242 between 2002 and 2012 (NEOREP Cohort, CNIL declaration number 1547270). We included unilateral,
243 non-recurrent, non-inflammatory, non-metastatic tumors, and excluded T4 tumors. This study was
244 conducted in accordance with institutional and ethical rules regarding research on tissue specimens and
245 patients. Information on family history, clinical characteristics (age; menopausal status; body mass
246 index) and tumor characteristics (clinical tumor stage and grade; histology; clinical nodal status; ER, PR
247 and *HER2* status; BC subtype; mitotic index; Ki67; lymphovascular invasion) was retrieved from

248 electronic medical records. All the patients of the cohort received NAC, and additional treatments were
249 decided in accordance with national guidelines.

250

251 ***Tumor samples and pathological review***

252 In accordance with the guidelines used in France (Group for Evaluation of Prognostic Factors
253 using Immunohistochemistry in Breast Cancer⁶¹), cases were considered estrogen receptor (ER)-positive
254 or progesterone receptor (PR)-positive if at least 10% of the tumor cells expressed estrogen and/or
255 progesterone receptors (ER/PR). Endocrine therapy was prescribed when this threshold was exceeded.
256 *HER2*-negative status was defined as a score of 0 or 1+ for the tissue section stained by
257 immunohistochemistry (IHC). Tissue sections with scores of IHC 2+ or IHC 3+ were then analyzed by
258 fluorescence *in situ* hybridization (FISH) to confirm *HER2* positivity. BC tumors were classified into
259 subtypes (TNBC, *HER2*-positive, and luminal *HER2*-negative [referred to hereafter as “luminal”]). BC
260 subtypes were defined as follows: luminal, ER⁺ or PR⁺/*HER2*⁻; TNBC, ER⁻/PR⁻/*HER2*⁻; *HER2*-positive
261 BC, *HER2*⁺. Pretreatment core needle biopsy specimens and/or the corresponding post-NAC surgical
262 specimens were reviewed independently by breast disease experts for research purposes, to assess
263 residual cancer burden index, and the levels of tumor-infiltrating lymphocytes. The pathological reviews
264 of these specimens are described in detail elsewhere^{62, 20, 63}. Pathological complete response (pCR) was
265 defined as the absence of residual invasive cancer cells in the breast and axillary lymph nodes (ypT0/is
266 p/ypN0).

267 ***Survival endpoints***

268 Relapse-free survival (RFS) was defined as the time from surgery to death, loco-regional
269 recurrence or distant recurrence, whichever occurred first. Overall survival (OS) was defined as the time
270 from surgery to death. The date of last known contact was retained for patients for whom none of these
271 events were recorded. The cutoff date for survival analysis was March, 13th, 2019.

272

273 ***Variables of interest***

274 The care pathway of BC patients eligible for neoadjuvant chemotherapy can be summarized as
275 follows: *i*) pretreatment biopsy for BC diagnosis; *ii*) administration of chemotherapy as the first-line
276 treatment; *iii*) removal of the tumor by surgery; *iv*) histological analysis of the specimens obtained; *v*)
277 prescription of adjuvant treatments, if indicated (radiotherapy, hormonotherapy, chemotherapy); *vi*)
278 patients follow-up to monitor for relapse or death. We identified 94 clinically relevant variables from
279 clinical, radiological, pathological and outcome data, which we grouped into 14 categories (hospital,
280 history, co-medication, comorbidities, clinical baseline, baseline histology, pre-NAC pathology,
281 treatment response, surgery, treatment, changes during NAC, post-NAC pathology, delayed
282 relapse/survival, metastasis). For composite variables derived from raw variables (e.g. BC subtype,
283 constructed from a combination of ER status, PR status, *HER2* status), both derived and raw variables
284 were represented on the network.

285

286 ***MIIC algorithm***

287 The functioning of the algorithm has been described in detail elsewhere^{5,7}. Briefly, starting from
288 a fully connected network, the MIIC algorithm first removes dispensable edges by iteratively subtracting
289 the most significant information contributions from indirect paths between each pair of variables. The
290 remaining edges, the underlying effect of which cannot be explained by indirect paths, are then oriented
291 based on the causality signature in the data, corresponding to the simultaneous head-to-head orientations
292 of so-called “v-structures”, $X \rightarrow Z \leftarrow Y$. In principle, propagation of v-structure orientations to
293 downstream edges can also be implemented to fulfill underlying model class assumptions^{64, 65} but are
294 not applied on the NEOREP clinical network to ensure that MIIC algorithmic decisions are only based
295 on information actually contained in the data.

296 Each edge corresponds to a “direct” association between two variables, that is, a statistical
297 association that cannot be entirely explained by indirect effects involving other variables. Red and blue
298 edges correspond to positive and negative (*i.e.* anti-correlated) associations, respectively. Four types of
299 edge orientations are distinguished by the MIIC online server: *i*) directed edges with a gray arrowhead

300 represent inferred causal relationships; *ii*) bidirected edges (drawn with dashed lines) reflect the presence
301 of a latent common cause (L) unobserved in the available dataset, *i.e.* $X \leftarrow (L) \rightarrow Y$; *iii*) directed edges
302 with a colored (red or blue) arrowhead are consistent with either a causal or a latent common cause
303 relationship; and *iv*) undirected edges, whose orientation if it exists, cannot be inferred from non-
304 perturbative data. The original algorithm was restricted to categorical variables⁵, but MIIC has recently
305 been extended to include continuous variables, the values of which are partitioned into optimal bins,
306 maximizing mutual information with another (continuous or categorical) variable of interest, while
307 preventing the overfitting of datasets of finite size due to the use of too many bins⁷. In particular, each
308 continuous variable may have different information-maximizing partitions depending on the associated
309 variable of interest. For instance, MIIC finds three maximally informative bins for the residual cancer
310 burden (RCB) score in association with patient survival status (**Fig. S1A**), whereas eight RCB bins are
311 required to estimate its mutual information with post-NAC cellularity correctly (**Fig. S1B**).

312

313 *MIIC online server*

314 The MIIC online server is freely accessible at <https://miic.curie.fr> and can be used with the
315 Google Chrome, Mozilla Firefox, Edge, and Safari browsers. The user guide summarizing the main
316 steps for running the MIIC algorithm is accessible at https://miic.curie.fr/user_guide.php, and an online
317 video tutorial is available at: <https://miic.curie.fr/tutorial.php>. The workbench is available from
318 <https://miic.curie.fr/workbench.php>. As input data, the user can upload a dataset formatted as a table
319 with commas, semicolons, tabs, pipes or colons, as separators, without row names. Each variable can be
320 either categorical or quantitative (discrete or continuous). Variables can be grouped into families,
321 identified with different colors on the network. Missing values are allowed in the dataset and their
322 possible statistical biases are taken into account by MIIC⁷. They should be indicated as “NA” in the
323 dataset table. Once the dataset has been prepared, the user runs the algorithm, and an e-mail is sent when
324 the job is completed.

325

326 *MIIC output*

327

328 The MIIC online server generates a visualization of the global network of the dataset. An example based
329 on the NEOREP dataset is displayed in Fig. 1, and is accessible as an interactive network at
330 https://miic.curie.fr/job_results_showcase.php?id=NEOREP.

331

332 *Interactive exploration of the network*

333 The distributions and neighborhoods of each node and edge of the inferred network can be explored
334 through an interactive interface, through the mouse-over right- or left-click buttons on the browser page,
335 as detailed in the online tutorials. Briefly, any variable can be highlighted by clicking on the network or
336 through the “Search” toolbox (Fig. S2A). The corresponding plots can be downloaded as .png or .svg
337 images. Each node can be explored individually in terms of counts (categorical variables, Fig. S2B-C)
338 or distribution (continuous variables Fig. S2D-E). Each edge can be explored by a right click and the
339 choice of “plot join distribution” or “plot discretization”. The resulting plots are (i) proportion plots,
340 with the edge representing the total association between two categorical variables (Fig. S2F); (ii)
341 distribution histograms (Fig. S2G) or boxplots (Fig. S2H), in which the edge represents the total
342 association between a categorical and a continuous variable or (iii) scatter plots (Fig. S2I), in which the
343 edge represents the total association between two continuous variables. Additional options include
344 inverting the x and y axes, the choice of frequency or absolute counts, or NA removal (proportion plots),
345 and faceting or superimposing the variables (distribution histograms). All the figures presented here
346 were generated with the MIIC online interactive visualization tool.

347

348 **Data availability statement:** All images and the associated network are publicly available at:
349 https://miic.curie.fr/job_results_showcase.php?id=NEOREP. Data corresponding to the NEOREP
350 cohort study will be available upon reasonable request.

351

352 **Acknowledgements:** NS acknowledges support from Sorbonne University (ATER), VC from ARC
353 foundation and HI from ITMO Cancer, Institut Curie and CNRS

354 **Author contributions:** NS, VC and HI designed and implemented the machine learning and
355 interactive exploration tools; ASH, LD, ML performed research; ASH and FR verified the data; NS,
356 VC, HI, LD, FR and ASH contributed to data analysis; ASH and FR contributed to expert review. All
357 authors contributed to data interpretation. LD, ASH, NS and HI wrote the paper. NS worked on the
358 paper while being affiliated at b. and c. All authors had full access to all the data in the study and had
359 final responsibility for the decision to submit for publication. ¹N.S., ²A-S.H. and ³V.C. **contributed**
360 **equally** to this work.

361

362 **Competing interests:** The authors declare no competing interests.

363

364 **References**

- 365 1. Bärtschi M: Health Data Visualization-A review * Seminar Collaborative Data Visualization, in
366 2015
- 367 2. Luo J, Wu M, Gopukumar D, et al: Big Data Application in Biomedical Research and Health Care:
368 A Literature Review. *Biomed Inform Insights* 8:1–10, 2016
- 369 3. Ola O, Sedig K: Beyond simple charts: Design of visualizations for big health data [Internet].
370 *Online J Public Health Inform* 8, 2016[cited 2019 Aug 14] Available from:
371 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5302463/>
- 372 4. Shneiderman B, Plaisant C, Hesse BW: Improving Healthcare with Interactive Visualization.
373 *Computer* 46:58–66, 2013
- 374 5. Verny L, Sella N, Affeldt S, et al: Learning causal networks with latent variables from multivariate
375 information in genomic data. *PLoS Comput Biol* 13:e1005662, 2017
- 376 6. Sella N, Verny L, Uguzzoni G, et al: MIIC online: a web server to reconstruct causal or non-causal
377 networks from non-perturbative data. *Bioinformatics* 34:2311–2313, 2018
- 378 7. Cabeli V, Verny L, Sella N, et al: Learning clinical networks from medical records based on
379 information estimates in mixed-type data [Internet]. *PLoS Comput Biol* 16, 2020[cited 2021 Feb 4]
380 Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7259796/>
- 381 8. Brandão M, Reyal F, Hamy A-S, et al: Neoadjuvant treatment for intermediate/high-risk HER2-
382 positive and triple-negative breast cancers: no longer an “option” but an ethical obligation. *ESMO*
383 *Open* 4:e000515, 2019
- 384 9. Reyal F, Hamy AS, Piccart MJ: Neoadjuvant treatment: the future of patients with breast cancer.
385 *ESMO Open* 3:e000371, 2018
- 386 10. Burstein HJ, Temin S, Anderson H, et al: Adjuvant Endocrine Therapy for Women With Hormone
387 Receptor–Positive Breast Cancer: American Society of Clinical Oncology Clinical Practice Guideline
388 Focused Update. *J Clin Oncol* 32:2255–2269, 2014
- 389 11. Wilson FR, Coombes ME, Wylie Q, et al: Herceptin® (trastuzumab) in HER2-positive early breast
14

390 cancer: protocol for a systematic review and cumulative network meta-analysis. *Syst Rev* 6:196, 2017
391 **12.** Aaldriks AA, Maartense E, Nortier HJWR, et al: Prognostic factors for the feasibility of
392 chemotherapy and the Geriatric Prognostic Index (GPI) as risk profile for mortality before
393 chemotherapy in the elderly. *Acta Oncol* 55:15–23, 2016
394 **13.** van Leeuwen RWF, Swart EL, Boven E, et al: Potential drug interactions in cancer therapy: a
395 prevalence study using an advanced screening method. *Ann Oncol* 22:2334–2341, 2011
396 **14.** Popa MA, Wallace KJ, Brunello A, et al: POTENTIAL DRUG INTERACTIONS AND
397 CHEMOTOXICITY IN OLDER PATIENTS WITH CANCER RECEIVING CHEMOTHERAPY. *J*
398 *Geriatr Oncol* 5:307–314, 2014
399 **15.** Wang N, Qin MZ, Cui J: Lipid profile comparison between pre- and post-menopausal women.
400 *Zhonghua Xin Xue Guan Bing Za Zhi* 44:799–804, 2016
401 **16.** Veronesi U, Paganelli G, Viale G, et al: Sentinel-lymph-node biopsy as a staging procedure in
402 breast cancer: update of a randomised controlled study. *Lancet Oncol* 7:983–990, 2006
403 **17.** Weidner N, Moore DH, Vartanian R: Correlation of Ki-67 antigen expression with mitotic figure
404 index and tumor grade in breast carcinomas using the novel “paraffin”-reactive MIB1 antibody. *Hum*
405 *Pathol* 25:337–342, 1994
406 **18.** Cortadellas T, Argacha P, Acosta J, et al: Estimation of tumor size in breast cancer comparing
407 clinical examination, mammography, ultrasound and MRI—correlation with the pathological analysis
408 of the surgical specimen. *Gland Surg* 6:330–335, 2017
409 **19.** Berg WA, Gutierrez L, NessAiver MS, et al: Diagnostic accuracy of mammography, clinical
410 examination, US, and MR imaging in preoperative assessment of breast cancer. *Radiology* 233:830–
411 849, 2004
412 **20.** Hamy A-S, Lam G-T, Laas E, et al: Lymphovascular invasion after neoadjuvant chemotherapy is
413 strongly associated with poor prognosis in breast carcinoma. *Breast Cancer Res Treat* 169:295–304,
414 2018
415 **21.** Meyers MO, Klauber-Demore N, Ollila DW, et al: Impact of breast cancer molecular subtypes on
416 locoregional recurrence in patients treated with neoadjuvant chemotherapy for locally advanced breast
417 cancer. *Ann Surg Oncol* 18:2851–2857, 2011
418 **22.** Lowery AJ, Kell MR, Glynn RW, et al: Locoregional recurrence after breast cancer surgery: a
419 systematic review by receptor phenotype. *Breast Cancer Res Treat* 133:831–841, 2012
420 **23.** Waljee JF, Hu ES, Newman LA, et al: Predictors of re-excision among women undergoing breast-
421 conserving surgery for cancer. *Ann Surg Oncol* 15:1297–1303, 2008
422 **24.** Truin W, Vugts G, Roumen RMH, et al: Differences in Response and Surgical Management with
423 Neoadjuvant Chemotherapy in Invasive Lobular Versus Ductal Breast Cancer. *Ann Surg Oncol*
424 23:51–57, 2016
425 **25.** Munhoz AM, Montag E, Gemperli R: Oncoplastic breast surgery: indications, techniques and
426 perspectives. *Gland Surg* 2:143–157, 2013
427 **26.** Buchholz TA: Radiation Therapy for Early-Stage Breast Cancer after Breast-Conserving Surgery.
428 *New England Journal of Medicine* 360:63–70, 2009
429 **27.** Carlson RW, Allred DC, Anderson BO, et al: Invasive breast cancer. *J Natl Compr Canc Netw*
430 9:136–222, 2011
431 **28.** Eifel P, Axelson JA, Costa J, et al: National Institutes of Health Consensus Development
432 Conference Statement: adjuvant therapy for breast cancer, November 1-3, 2000. *J Natl Cancer Inst*
433 93:979–989, 2001
434 **29.** Halberg FE, Shank BM, Haffty BG, et al: Conservative surgery and radiation in the treatment of
435 stage I and II carcinoma of the breast. American College of Radiology. ACR Appropriateness Criteria.
436 *Radiology* 215 Suppl:1193–1205, 2000
437 **30.** Masuda N, Lee S-J, Ohtani S, et al: Adjuvant Capecitabine for Breast Cancer after Preoperative
438 Chemotherapy. *N Engl J Med* 376:2147–2159, 2017
439 **31.** Dent R, Valentini A, Hanna W, et al: Factors associated with breast cancer mortality after local
440 recurrence. *Curr Oncol* 21:e418–e425, 2014
441 **32.** Whelan T, Clark R, Roberts R, et al: Ipsilateral breast tumor recurrence postlumpectomy is
442 predictive of subsequent mortality: results from a randomized trial. Investigators of the Ontario

- 443 Clinical Oncology Group. *Int J Radiat Oncol Biol Phys* 30:11–16, 1994
- 444 **33.** Kurtz JM, Spitalier JM, Amalric R, et al: The prognostic significance of late local recurrence after
445 breast-conserving therapy. *Int J Radiat Oncol Biol Phys* 18:87–93, 1990
- 446 **34.** Sopik V, Nofech-Mozes S, Sun P, et al: The relationship between local recurrence and death in
447 early-stage breast cancer. *Breast Cancer Res Treat* 155:175–185, 2016
- 448 **35.** Witteveen A, Kwast ABG, Sonke GS, et al: Survival after Locoregional Recurrence or Second
449 Primary Breast Cancer: Impact of the Disease-Free Interval. *PLOS ONE* 10:e0120832, 2015
- 450 **36.** Hudis CA, Barlow WE, Costantino JP, et al: Proposal for standardized definitions for efficacy end
451 points in adjuvant breast cancer trials: the STEEP system. *J Clin Oncol* 25:2127–2132, 2007
- 452 **37.** Voduc KD, Cheang MCU, Tyldesley S, et al: Breast Cancer Subtypes and the Risk of Local and
453 Regional Relapse. *JCO* 28:1684–1691, 2010
- 454 **38.** Wu X, Baig A, Kasymjanova G, et al: Pattern of Local Recurrence and Distant Metastasis in
455 Breast Cancer By Molecular Subtype [Internet]. *Cureus* 8, 2016[cited 2021 Feb 4] Available from:
456 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5222631/>
- 457 **39.** Snell CE, Gough M, Middleton K, et al: Absent progesterone receptor expression in the lymph
458 node metastases of ER-positive, HER2-negative breast cancer is associated with relapse on tamoxifen.
459 *J Clin Pathol* 70:954–960, 2017
- 460 **40.** Nishimura R, Osako T, Okumura Y, et al: Changes in the ER, PgR, HER2, p53 and Ki-67
461 biological markers between primary and recurrent breast cancer: discordance rates and prognosis.
462 *World J Surg Oncol* 9:131, 2011
- 463 **41.** Nishimura R, Osako T, Nishiyama Y, et al: Evaluation of factors related to late recurrence--later
464 than 10 years after the initial treatment--in primary breast cancer. *Oncology* 85:100–110, 2013
- 465 **42.** Darlix A, Griguolo G, Thezenas S, et al: Hormone receptors status: a strong determinant of the
466 kinetics of brain metastases occurrence compared with HER2 status in breast cancer. *J Neurooncol*
467 138:369–382, 2018
- 468 **43.** Zhou L, Zhou W, Zhang H, et al: Progesterone suppresses triple-negative breast cancer growth and
469 metastasis to the brain via membrane progesterone receptor α . *Int J Mol Med* 40:755–761, 2017
- 470 **44.** Ahern TP, Pedersen L, Tarp M, et al: Statin prescriptions and breast cancer recurrence risk: a
471 Danish nationwide prospective cohort study. *J Natl Cancer Inst* 103:1461–1468, 2011
- 472 **45.** Kwan ML, Habel LA, Slattery ML, et al: NSAIDs and Breast Cancer Recurrence in a Prospective
473 Cohort Study. *Cancer Causes Control* 18:613–620, 2007
- 474 **46.** Powe DG, Voss MJ, Zänker KS, et al: Beta-Blocker Drug Therapy Reduces Secondary Cancer
475 Formation in Breast Cancer and Improves Cancer Specific Survival. *Oncotarget* 1:628–638, 2010
- 476 **47.** Hamy A-S, Derosa L, Valdelièvre C, et al: Comedications influence immune infiltration and
477 pathological response to neoadjuvant chemotherapy in breast cancer. *OncoImmunology* 9:1677427,
478 2020
- 479 **48.** Farrugia DJ, Landmann A, Diego E, et al: Mitotic index to predict breast cancer recurrence after
480 neoadjuvant systemic therapy. *JCO* 34:e23265–e23265, 2016
- 481 **49.** Pattali S, Harding N, Visotcky A, et al: Value of mitotic index in residual tumors following
482 neoadjuvant therapy for breast cancer: Single institution experience. *JCO* 34:548–548, 2016
- 483 **50.** Symmans WF, Peintinger F, Hatzis C, et al: Measurement of Residual Breast Cancer Burden to
484 Predict Survival After Neoadjuvant Chemotherapy. *Journal of Clinical Oncology* 25:4414–4422, 2007
- 485 **51.** Duggan MA, Anderson WF, Altekruse S, et al: The Surveillance, Epidemiology and End Results
486 (SEER) Program and Pathology: Towards Strengthening the Critical Relationship. *Am J Surg Pathol*
487 40:e94–e102, 2016
- 488 **52.** James B. Yu MD: NCI SEER Public-Use Data: Applications and Limitations in Oncology
489 Research [Internet]. *Cancer Network* , 2009[cited 2019 Aug 27] Available from:
490 [https://www.cancernetwork.com/oncology-journal/nci-seer-public-use-data-applications-and-](https://www.cancernetwork.com/oncology-journal/nci-seer-public-use-data-applications-and-limitations-oncology-research)
491 [limitations-oncology-research](https://www.cancernetwork.com/oncology-journal/nci-seer-public-use-data-applications-and-limitations-oncology-research)
- 492 **53.** Boffa DJ, Rosen JE, Mallin K, et al: Using the National Cancer Database for Outcomes Research:
493 A Review. *JAMA Oncol* 3:1722–1728, 2017
- 494 **54.** Bezin J, Duong M, Lassalle R, et al: The national healthcare system claims databases in France,
495 SNIIRAM and EGB: Powerful tools for pharmacoepidemiology. *Pharmacoepidemiol Drug Saf*

496 26:954–962, 2017

497 **55.** Tuppin P, Rudant J, Constantinou P, et al: Value of a national administrative database to guide
498 public decisions: From the système national d’information interrégimes de l’Assurance Maladie
499 (SNIIRAM) to the système national des données de santé (SNDS) in France [Internet].
500 /data/revues/03987620/v65sS4/S0398762017304315/ , 2017[cited 2019 Aug 13] Available from:
501 <https://www.em-consulte.com/en/article/1140905>

502 **56.** Margolis R, Derr L, Dunn M, et al: The National Institutes of Health’s Big Data to Knowledge
503 (BD2K) initiative: capitalizing on biomedical big data. *J Am Med Inform Assoc* 21:957–958, 2014

504 **57.** Keim D, Andrienko G, Fekete J-D, et al: Visual Analytics: Definition, Process, and Challenges
505 [Internet], in Kerren A, Stasko JT, Fekete J-D, et al (eds): *Information Visualization*. Berlin,
506 Heidelberg, Springer Berlin Heidelberg, 2008, pp 154–175[cited 2019 Aug 14] Available from:
507 http://link.springer.com/10.1007/978-3-540-70956-5_7

508 **58.** Happe A, Drezen E: A visual approach of care pathways from the French nationwide SNDS
509 database - from population to individual records: the ePEPS toolbox [Internet], 2018[cited 2019 Aug
510 18] Available from: <https://hal-univ-rennes1.archives-ouvertes.fr/hal-01697626>

511 **59.** Zhao Y, Parvinzamid F, Wei H, et al: Visual Analytics for Health Monitoring and Risk
512 Management in CARRE. *E-Learning and Games; 10th International Conference, Edutainment 2016*,
513 Hangzhou, China, April 14-16, 2016, Revised Selected Papers 9654:380–391, 2016

514 **60.** Ledesma A, Al-Musawi M, Nieminen H: Health figures: an open source JavaScript library for
515 health data visualization [Internet]. *BMC Med Inform Decis Mak* 16, 2016[cited 2019 Aug 14]
516 Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4802654/>

517 **61.** [Recommendations for the immunohistochemistry of the hormonal receptors on paraffin sections
518 in breast cancer. Update 1999. Group for Evaluation of Prognostic Factors using
519 Immunohistochemistry in Breast Cancer (GEFPICS-FNCLCC)]. *Ann Pathol* 19:336–343, 1999

520 **62.** Hamy A-S, Pierga J-Y, Sabaila A, et al: Stromal lymphocyte infiltration after neoadjuvant
521 chemotherapy is associated with aggressive residual disease and lower disease-free survival in HER2-
522 positive breast cancer. *Ann Oncol* 28:2233–2240, 2017

523 **63.** Hamy-Petit A-S, Belin L, Bonsang-Kitzis H, et al: Pathological complete response and prognosis
524 after neoadjuvant chemotherapy for HER2-positive breast cancers before and after trastuzumab era:
525 results from a real-life cohort. *Br J Cancer* 114:44–52, 2016

526 **64.** Affeldt S, Isambert H: Robust reconstruction of causal graphical models based on conditional 2-
527 point and 3-point information, in *Proceedings of the 31th conference on Uncertainty in Artificial
528 Intelligence (UAI)*, 2015

529 **65.** Affeldt S, Verny L, Isambert H: 3off2: A network reconstruction algorithm based on 2-point and
530 3-point information statistics. *BMC Bioinformatics* 17 Suppl 2:12, 2016

531

532

533 **Figures Legends**

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

Figure 1: MIIC global network for the NEOREP breast cancer cohort. Each node corresponds to a variable of the dataset, with circles indicating continuous variables and squares indicating categorical variables. The colors define a category of variables, as detailed under the figure. Each edge corresponds to a “direct” association between two variables with different types of orientation described in Methods. *BC= breast cancer, BMI= body mass index, DCIS=ductal carcinoma in situ, ER= estrogen receptor status, LVI= lymphovascular invasion, NAC= neoadjuvant chemotherapy, CNS= central nervous system, pCR= pathological complete response, PR=progesterone receptor status, RCB= residual cancer burden, TILs= tumor-infiltrating lymphocytes.* Blue edges indicate negative partial correlations, red edges indicate positive partial correlations.

Figure 2. The MIIC interactive online interface identifies inherent associations between variables.

a) NAC type is directly correlated with NAC duration. NAC=neoadjuvant chemotherapy b) Distribution of neoadjuvant chemotherapy (NAC) duration (in days) according to the NAC regimen administered: anthracyclines (AC), taxanes or sequential AC-taxanes c) The number of axillary nodes in the histological specimen depends on the type of axillary surgery performed d) Boxplot showing the number of axillary nodes removed according to the type of surgery performed: lymph node dissection (LND), sentinel lymph node biopsy (SLN) or both e) Network interactions of the RCB node with the five patterns making up the RCB score.

Figure 3. MIIC identifies differences in clinical practices between the two centers of the cohort

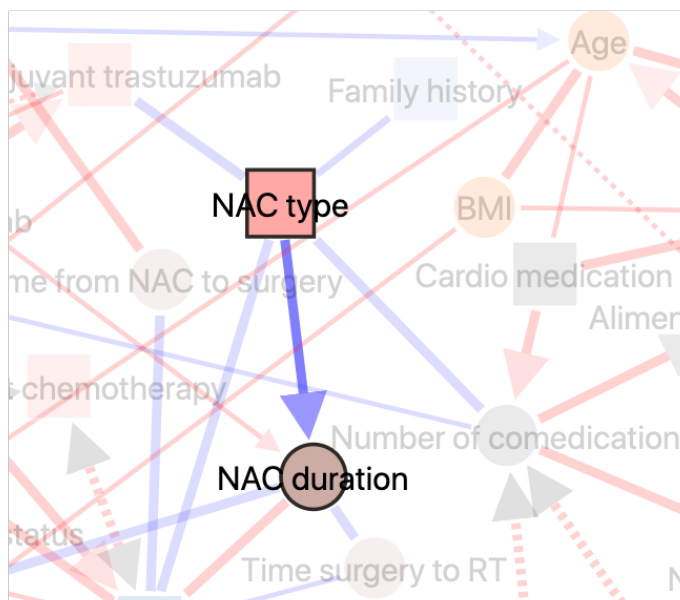
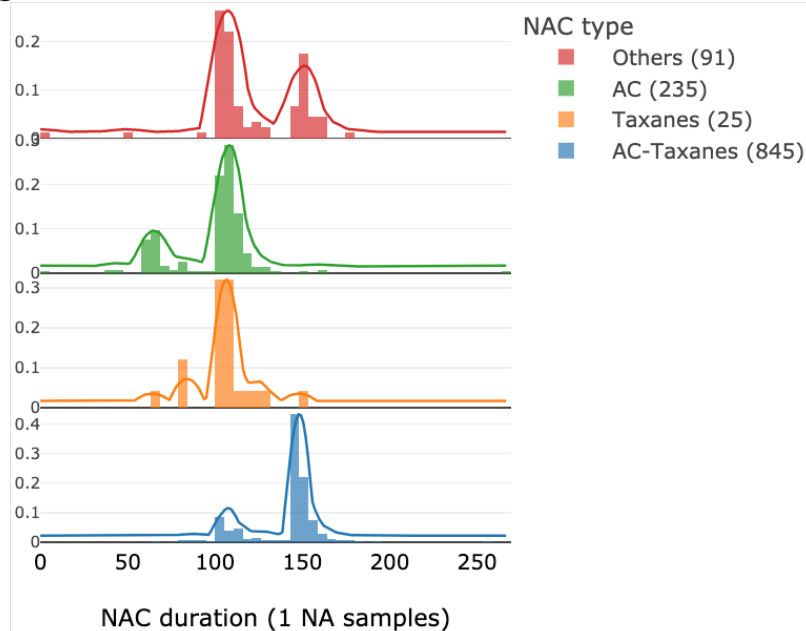
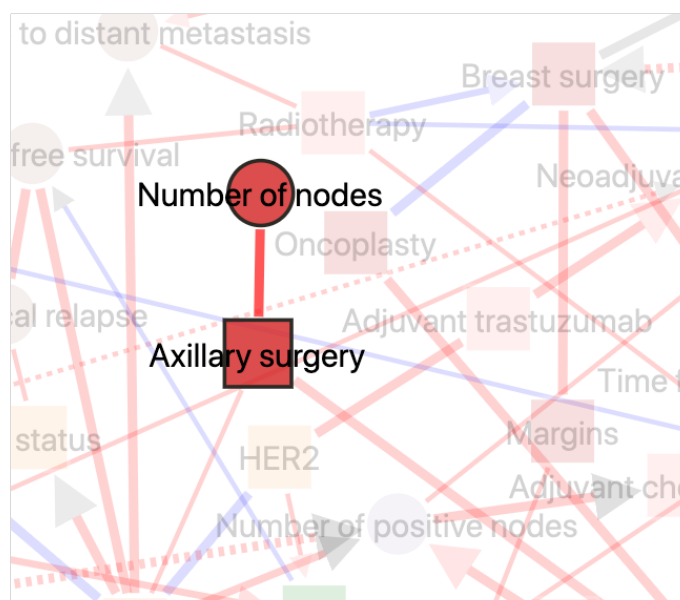
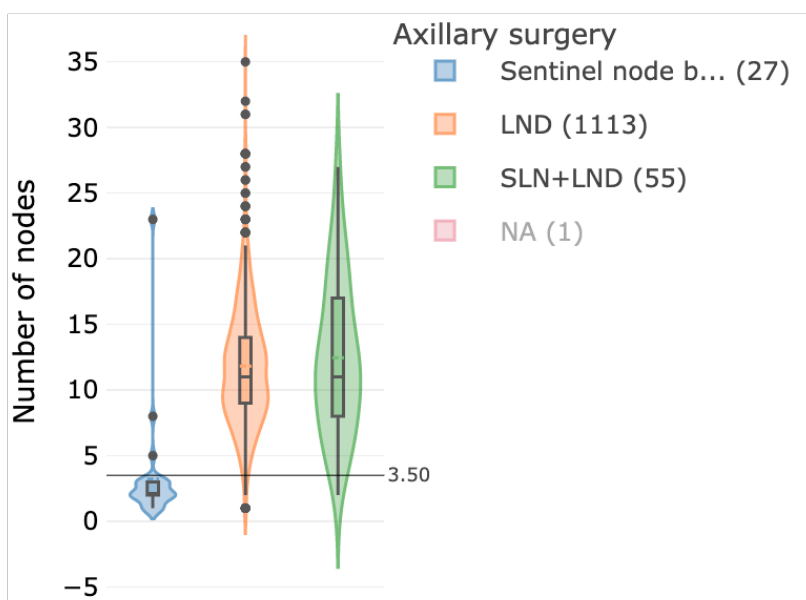
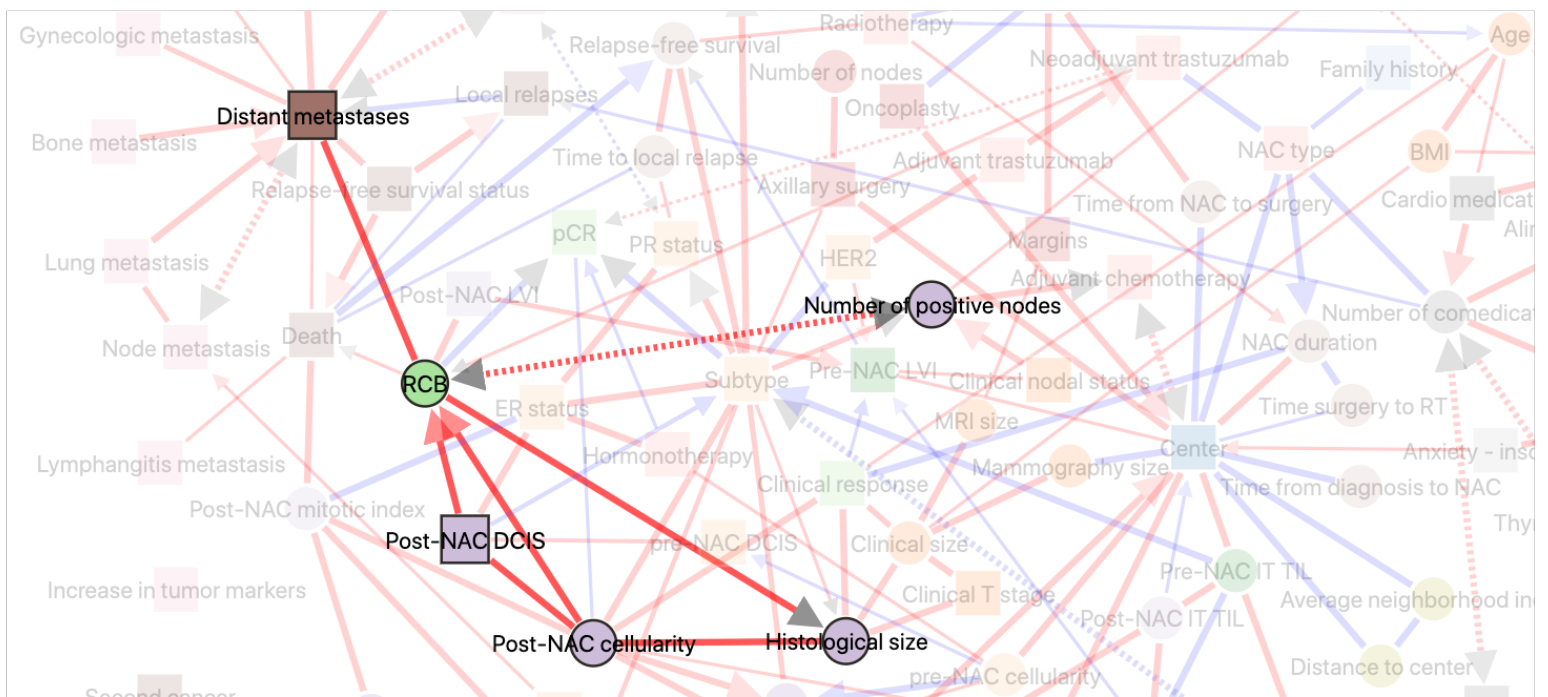
a) Network interactions around the node “center” of treatment. b) Proportion of patients undergoing oncoplastic surgery, according to treatment center: Paris or St Cloud c) Proportion of patients receiving adjuvant chemotherapy according to treatment center: Paris or St Cloud. d) Proportion of the various NAC regimens according to treatment center. e) Distribution plot for NAC duration in days, according to treatment center.

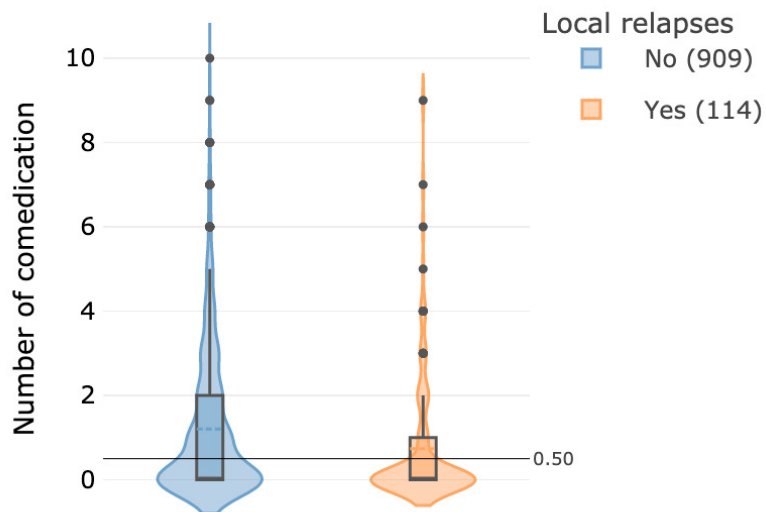
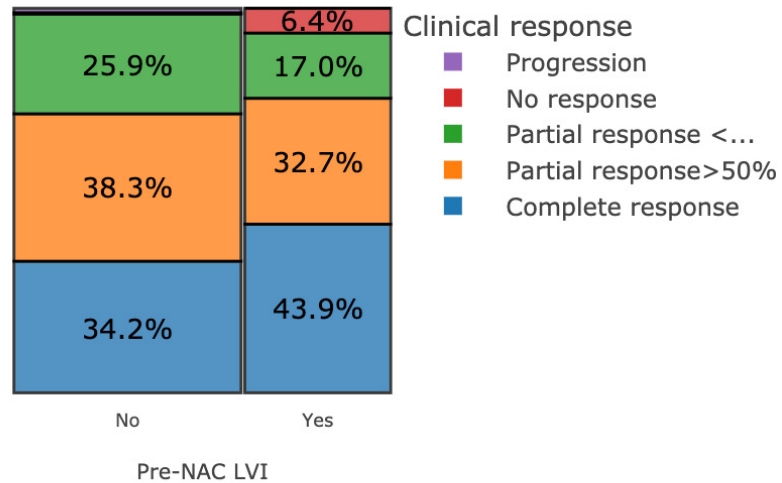
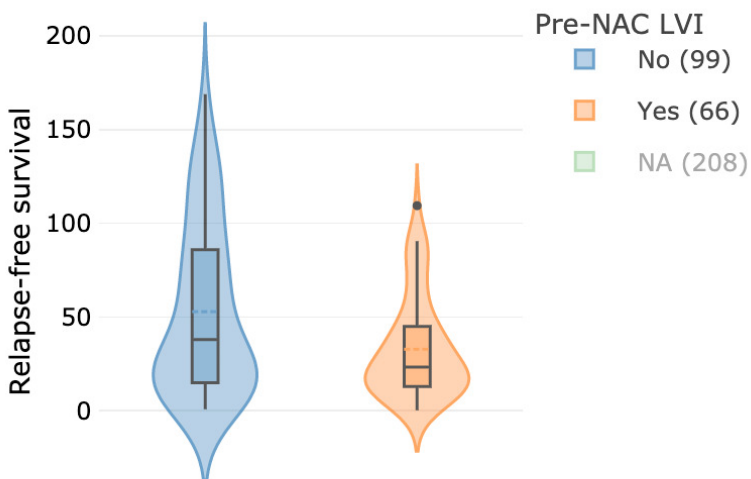
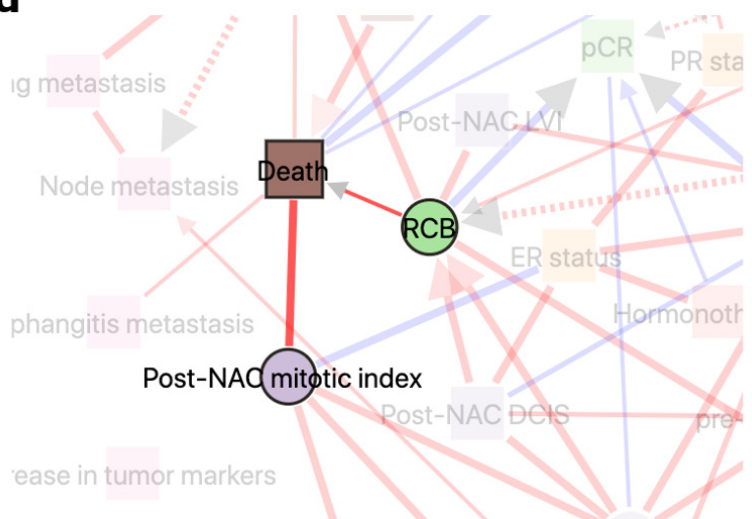
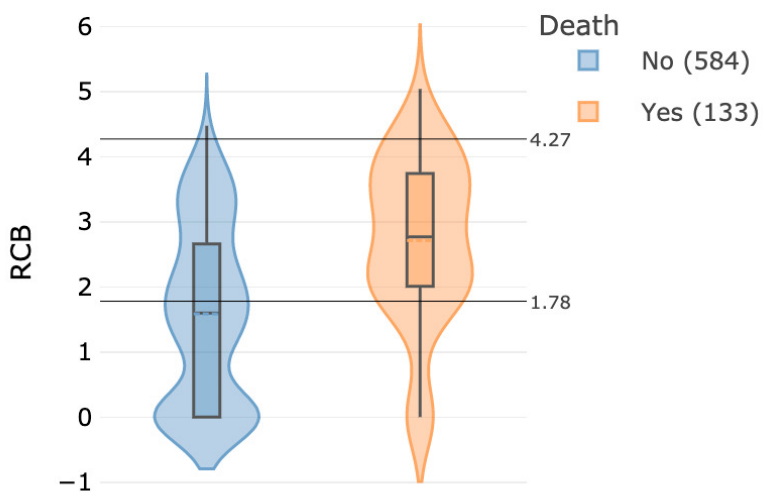
Figure 4. MIIC traces the natural course of the disease

a) Network interactions showing links between relapses, metastases and death in breast cancer. b) Proportion of distant metastases according to the occurrence or absence of local relapses. c) Proportion of deaths according to distant metastasis status. d) Distribution plot for relapse-free survival (in months) according to breast cancer subtype. e) Proportion plot displaying the relationship between central nervous system (CNS) metastasis and progesterone receptor (PR) status.

Figure 5 MIIC identifies factors likely to improve prediction or prognosis.

a) Network interaction displaying the link between local relapse occurrence and the number of drugs taken (comedication). b) Proportion plot showing the percentage of different clinical responses according to the presence or absence of pre-NAC lymphovascular invasion. c) Boxplot of relapse-free survival according to the presence or absence of pre-NAC lymphovascular invasion. d) Network interaction displaying the link between death, RCB and post-NAC mitotic index. e) Boxplot of RCB values according to vital status. f) Boxplot of post-NAC mitotic index according to vital status.

a**b****c****d****e**

a**b****c****d****e****f**