

Community Science and Reaching the Promise of Big Data in Health Care

Charles Mayo, PhD

University of Michigan

cmayo@med.umich.edu

Disclosures: None

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/mp.13140](https://doi.org/10.1002/mp.13140)

This article is protected by copyright. All rights reserved

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

Article Type: Editorial

“Alone we can do so little; together we can do so much” - Helen Keller

At its core, the concept of “big data” in health care embraces the promise of creating transcendent knowledge generation systems using the power of information gathered from routine processes for all patients. By harnessing large-scale aggregation of information generated during routine healthcare delivery to the speed and-capacity of machine learning (ML) and artificial intelligence (AI) algorithms, “big data” pioneers believe we can reduce research costs, deepen our understanding of factors affecting patient outcomes and improve patient care. Ability of ML and AI to identify and characterize interactions among sets of variables and cohorts of patients, much larger than humans are capable of analyzing, should deepen our understanding key factors affecting outcomes and quality of care. These are early days; with substantial amounts of foundational work needed to reach that promise. Part of that foundation includes improving standardizations for quantifying data elements and building systems to increase the volume of quality data that may be consumed by these data hungry ML and AI algorithms.

By analogy, consider our love of fast, powerful and convenient cars. Without the foundational work of constructing roads, bridges, fuel stations and other supports, connections and logistics plus development of science and engineering underlying design, the reality of driving would not be possible. A myriad of standardizations (e.g. diameter of nozzle on fuel pump at gas station, traffic laws, road design, regulatory standards) enable us to just focus on driving, without having to also grapple with endless variations in key details. These elements did not emerge quickly and fully formed from the minds of a handful of people. Instead, they evolved, gradually out of the combined trial and error iterations of communities of enthusiasts to find a common solution.

Similarly making the promise of “big data” a practical, routine part of our clinical reality will be an outgrowth of what we are able to do together as a community to build needed core concepts (e.g. clinically linked measures of ML/AI algorithm reliability) and standardizations (e.g. nomenclatures, ontologies, toxicity measures, disease site status/recurrence categorizations). Practice standardizations (e.g. how recurrence information is entered into a treatment note) enable our electronic systems to

33 make distinctions among data elements that can then be fed accurately, rapidly and in large volume to
34 learning algorithms. In our cars, standardizations let us take for granted the ability to drive up to any gas
35 station to fuel our travels. By contrast, lack of standardized categorizations and entry processes in our
36 clinics, means we cannot take for granted the ability to electronically extract accurate information on
37 treatment outcomes, treatment variables, and relevant patient host variables from available electronic
38 health records to fuel our treatment outcomes modeling.

39
40 Further, if we aspire to eventually understand global patterns of care and treatment outcomes for all
41 cancer patients treated by our healthcare systems, as opposed to outcomes of limited patient cohorts
42 accrued at a relatively small number of centers, then our communities and the foundational work
43 required of them, must expand beyond the scale of a few institutions. It requires that we move toward
44 community science, where collaborations spanning multiple institutions, clinic sizes and national
45 borders are recognized as key factors for success in supporting creation of practical enabling
46 standardizations, ontologies, algorithms and processes. While this principle is recognized in clinical trials,
47 the comparable funding, institutional and academic supports for the foundational “big data” efforts
48 often are not. Using our analogy, it is often easier to get support for designing a futuristic car than for
49 constructing roads and bridges.

50
51 One recent example of community science is the American Association of Physicists in Medicine’s
52 (AAPM’s) Task Group 263 (TG-263) on Radiation Oncology Nomenclature.[1] This task group worked
53 with a large and diverse group of 57 physicians, physicists, industry representatives and others, drawn
54 from large clinics and small, academic and non-academic centers, the AAPM, the American Society of
55 Therapeutic Radiation Oncology (ASTRO), European Society for Radiation Oncology (ESTRO), NRG and
56 IHE-RO and other stakeholder groups. The Task Group created and piloted a proposed set of
57 nomenclature standardization recommendations designed to improve the ability to electronically
58 extract and use large data sets of dosimetric data to support “big data” efforts. For example, when
59 analyzing the history of treatment plans at an institution one often finds dozens of character
60 combinations that are used to represent each organ at risk (e.g. left optic nerve). Once the
61 nomenclature is in place, with just one recommended naming for each structure, it is possible to
62 automate accurate, routine extraction of large volumes of dosimetric data on treatment planning
63 structures for analysis. An important lesson from that effort was the vital role that professional societies
64 play in supporting and endorsing these efforts.

65

66 This special issue of *Medical Physics* is another example of a community effort to bring the promise of
67 “big data” closer to reality. The first Practical Big Data Workshop (PBDW), held at the University of
68 Michigan in Ann Arbor in 2017, was an effort to promote coalescing the nascent community of builders
69 and users of “big data” systems for cancer care. Shared recognition of a common set of challenges and
70 need for consensus solutions gave rise to a set of papers to share these perspectives with the larger
71 community.

72

73 As healthcare builds slowly toward a reality where “big data” and analytics become more routine
74 elements of clinical practice, there are significant implications for the training and credentialing of
75 healthcare professionals. A lesson from the PBDW is that physicians, physicists and others must blend
76 their Radiation Oncology domain knowledge with skill sets from other domains (e.g. informatics,
77 application development, machine learning, ethics, genomics, radiomics, etc.). New combinations of skill
78 sets, e.g. physician-ethicists, physician-informaticists, physicist-data scientists, and physicist-database
79 designers, played vital roles at the meeting in identifying challenges, formulating solutions and
80 effectively communicating these challenges and solutions to the wider community.

81

82 Realizing the promise of “big data” in health care will be most effectively approached if we work
83 together as communities transcending boundaries that now separate institutions, professional
84 identities, and differently structured clinical service lines. Embracing both the need for expanding the
85 range of skill sets outside our traditional health care training curricula and the importance of building
86 networks of collaborators will enable us to build the strong foundation needed for knowledge
87 generating systems to emerge.

88

89 The papers in this special issue span a wide range of subject areas encountered in meeting specific
90 challenges, solutions and collaborative efforts which are part of reaching the potential of “big data” in
91 radiation oncology.

92

93 In “Treatment Data and Technical Process Challenges for Practical Big Data Efforts in Radiation
94 Oncology”, Mayo et al address several factors affecting many key data elements. [2] These include: need
95 for process and system changes to improve quality and availability of key data elements and
96 relationships, access and extraction issues for obtaining data from various source systems used in

97 patient treatment, selection considerations for database technologies, review and comparison of clinical
98 data repositories, specific recommendations for workflows and standardizations, examination of next
99 steps needed to improve data availability. In addition, the appendix of this manuscript details a
100 translational research ontology that specifies core data elements and relationships important to a broad
101 range of patient quality improvement and translational research efforts. Their recommendations for
102 improving clinical process include: more complete and consistent utilization of diagnosis and staging
103 tools in radiation oncology information systems, implementation of TG-263 standardizations for
104 nomenclature, routine creation of as treated plan sums in treatment planning systems that reflect all
105 dose delivered in the treatment course and standardized entry of patient reported outcomes and
106 provider reported toxicities into the electronic record.

107
108 Matuszak et al focus on the efforts and challenges for aggregation of outcomes information in their
109 manuscript, "Performance/Outcomes Data and Physician Process Challenges for Practical Big Data
110 Efforts in Radiation Oncology". [3] Building from a detailed examination of the "big data" projects of 8
111 groups, they examine common issues affecting data availability, access, and quality. They provide
112 specific recommendations for improvements through standardized workflows and discuss need for
113 multi-institutional consensus based standards for classifying recurrence categorizations.

114
115 In "Genomics, Bio specimens and other Biological Data: Current status and future directions
116 " by Rosenstein et al, challenges for the use of genomic and bio-specimen data are examined. [4]
117 Acquisition and storage of this key data element is currently the exception. They examine state large
118 scale research efforts to using this data element, challenges for access and extraction, issues for
119 collection and curation and provide specific recommendations for standardizations aimed at reducing
120 barriers to more wide spread, routine use of this data to support modeling patient outcomes.
121 Recommendations include developing a standardized nomenclature to reduce variability in collecting
122 genomics and bio-specimen data, developing standardizations in through multi-institutional and vendor
123 collaborations to improve interoperability, increasing the frequency of multi-institutional data pooling,
124 and harmonizing approaches for encapsulating this information in the EHR.

125
126 Mackie et al deal with challenges for "big data "in aggregation of imaging information, radiomics
127 measures and analysis of quantitative images to find biomarkers of disease in "Opportunities and
128 Challenges to Utilization of Quantitative Imaging: Report of the AAPM Practical Big Data Workshop".[5]

129 They address challenges in information curation that stand as obstacles to medicine transforming itself
130 into a “knowledge-based” discipline, carefully referencing impacts on clinical trials and NCI funded
131 imaging consortia. Highlights from their recommendations include need: for Radiology practices
132 embrace the needs of Oncology for more detailed quantitative imaging features, to include more
133 quantitative measurement data in images, to improve standardized radiology and oncology workflows,
134 to add more quantitative information on image features as part of routine practice, to improve
135 quantitative imaging reproducibility, accuracy and curation, and to examine approaches to regulation of
136 imaging biomarkers

137
138 In “Machine Learning and Modeling: Data, Validation, Communication Challenges”, El Naqa et al
139 highlight the potential of ML and AI for clinical advancement, while also addressing pitfalls when
140 applying these powerful analytic tools.[6] They discuss common issues requiring careful consideration
141 including proper use of analysis metrics, sufficient volume and quality of data in training sets, parsimony
142 and generalizability of models, quality assurance and clinical interpretability of results.
143 Recommendations include, establishing standardized clinically relevant objective criteria for evaluating
144 ML results, constructing publically available benchmark data sets to validate and cross check models,
145 using resampling techniques to estimate model performance, and benchmarking changes in predictive
146 performance of ML models using new biomarkers against with comparison to standard clinical factors.

147
148 The move to construction of learning health systems requires careful consideration of ethical obligations
149 to patients, construction of informed consent and addressing inconsistencies and variable
150 interpretations of the regulatory environment. Spector-Bagdady and Jagsi provided much need guidance
151 and perspective for addressing these challenges in “Big Data, Ethics, and Regulations: Implications for
152 Consent in the Learning Health System”. [7]

153
154 Traverso et al discuss their extensive experience with multi-institutional data sharing practices in “The
155 Radiation Oncology Ontology (ROO): publishing linked data in radiation oncology using Semantic Web
156 and Ontology techniques”. [8] Use of standardizations and design of scalable “big data” systems are
157 important principles for making data sets Findable, Accessible, Interoperable and Reusable (FAIR). They
158 discuss their use of the ROO with semantic web technologies to meet these goals.

159

160 United States based clinicians and researchers may be unfamiliar with the extensive efforts in Canada to
161 improve the landscape for “big data”, as part of improving safety and practice quality initiatives. In
162 “Improving Patient Outcomes and Radiotherapy Systems: A Pan-Canadian Approach to Patient Reported
163 Outcome Use”, Caissie et al provide a review of these efforts. [9] Among them, the Canadian Partnership
164 for Radiation Therapy (CPQR) combines the work of several groups including Association of Radiation
165 Oncology (CARO), Canadian Organization of Medical Physicists (COMP), and the Canadian Association of
166 Medical Radiation Technologists (CAMRT). The work of CPQR in promulgating standardizations (e.g. TG-
167 263) and key quality indicators is discussed. In addition, their innovative work developing standardized
168 approaches to administration and use of patient reported outcomes (PROs) across Canada is presented.

169
170 In “Practical data collection and extraction for Big Data applications in radiotherapy”, McNutt et al
171 discuss their experience with overcoming practical challenges to capture of high quality treatment and
172 outcomes data. Detailed examinations of factors affecting: clinician assessments, PROs, bio-specimen,
173 imaging, treatment and symptom management are discussed. They review approaches to technology
174 and clinical implementation they have used to address these challenges.

175
176 In “Perspectives on potential research benefits from big data efforts in Radiation Oncology”, Vikram
177 discusses several themes that frequently emerge in research efforts that may be positively affected by
178 these “big data” efforts. [11] Specific challenges facing radiation oncology and areas that “big data”
179 researchers should try to address are discussed.

180
181 Wei et al examine the impact of “big data” efforts and supporting standardizations on clinical trials in
182 “Implementation and enforcement of the standardization for radiotherapy with protocol guidelines,
183 libraries and software systems assure the clinical trial data quality”. [12] They share their in-depth
184 perspective on implementation details of management tools in the several network groups of the
185 National Clinical Trials Network (NCTN). From this perspective they underscore the significant overlap of
186 standardization recommendations highlighted throughout this special issue with NCTN objectives.

187
188 We hope that you will find the manuscripts in this special issue helpful in your personal journey and
189 entry into this growing community of practical big data in health care.

190

191 **References**

192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221

- 1) Mayo CS, Moran JM, Bosch W, Xiao Y, McNutt T, Popple R, Michalski J, Feng M, Marks LB, Fuller CD, Yorke E, Palta J, Gabriel PE, Molineu A, Matuszak MM, Covington E, Masi K, Richardson SL, Ritter T, Morgas T, Flampouri S, Santanam L, Moore JA, Purdie TG, Miller RC, Hurkmans C, Adams J, Jackie Wu QR, Fox CJ, Siochi RA, Brown NL, Verbakel W, Archambault Y, Chmura SJ, Dekker AL, Eagle DG, Fitzgerald TJ, Hong T, Kapoor R, Lansing B, Jolly S, Napolitano ME, Percy J, Rose MS, Siddiqui S, Schadt C, Simon WE, Straube WL, St James ST, Ulin K, Yom SS, Yock TI: American Association of Physicists in Medicine Task Group 263: Standardizing Nomenclatures in Radiation Oncology Int Journal Radiat Oncol Biol Phys 100(4): 1057-1066, 2018. PMID: 29485047
- 2) Mayo CS, Phillips M, McNutt T, Palta J, Dekker A, Miller R, Xiao Y, Moran JM, Matuszak MM, Gabriel P, Ahmet A, Prisciandaro J, Thor M, Phillips M, Dixit N, Popple R, Killoran J, Kaleba E, Kantor M, Ruan D, Kapoor R, Kessler M, Lawrence T: Treatment Data and technical Process Challenges for Practical Big Data Efforts in Radiation Oncology. Medical Physics (Accepted for PBDW2017 Issue): 2018. (In Press)
- 3) Martha Marie Matuszak, Clifton David Fuller, Torunn Yock, Clayton B Hess, Todd R McNutt, Shruti Jolly, Peter Gabriel, Charles Mayo, Maria Thor, Amanda Caissie, Arvind Rao, Dawn Owen, Wade P Smith, Jatinder R. Palta, Rishabh Kapoor, James Hayman, Mark Waddle, Barry Rosenstein, Robert Miller, Seungtaek Choi, Amy Moreno, Joseph Herman, Mary Feng: Performance/Outcomes Data and Physician Process Challenges for Practical Big Data Efforts in Radiation Oncology. Medical Physics (Accepted for PBDW2017 Issue): 2018. (In Press)
- 4) Rosenstein B, Rao A, Moran JM, Spratt DE, Mendonca M, Al-Lazikani B, Mayo CS, Speers C. Genomics, Bio specimens and other Biological Data: Current status and future directions. Medical Physics (Accepted for PBDW2017 Issue): 2018. (In Press)

- 222 5) Thomas R. Mackie, Edward F. Jackson¹, Maryellen Giger. Opportunities and Challenges to
223 Utilization of Quantitative Imaging Report of the AAPM Practical Big Data Workshop. Medical
224 Physics (Accepted for PBDW2017 Issue): 2018. (In Press)
225
226
- 227 6) El Naga I, Ruan D, Valdes G, Dekker A, McNutt T, Ge Y, Qing-Wu QR, Oh JH, Thor M, Smith W,
228 Rao A, Fuller C, Xiao Y, Manion F, Schipper M, Mayo CS, Moran JM, Ten Haken RK.: Machine
229 Learning and Modeling: Data, Validation, Communication Challenges. Medical Physics
230 (Accepted for PBDW2017 Issue): 2018. (In Press)
231
232
- 233 7) Kayte Spector-Bagdady, Reshma Jagsi. Big Data, Ethics, and Regulations: Implications for
234 Consent in the Learning Health System. Medical Physics. (Accepted for PBDW2017 Issue): 2018.
235 (In Press)
236
- 237 8) Alberto Traverso, Johan van Soest, *, Leonard Wee and Andre Dekker. The Radiation Oncology
238 Ontology (ROO): publishing linked data in radiation oncology using Semantic Web and
239 Ontology techniques. Medical Physics. (Accepted for PBDW2017 Issue): 2018. (In Press)
240
- 241 9) Amanda Caissie, Erika Brown, Rob Olson, Lisa Barbera, Carol-Anne Davis, Michael Brundage,
242 Michael Milosevic. Improving Patient Outcomes and Radiotherapy Systems: A Pan-Canadian
243 Approach to Patient Reported Outcome Use. Medical Physics. (Accepted for PBDW2017 Issue):
244 2018. (In Press)
245
- 246 10) Todd R. McNutt, Michael Bowers, Sierra Cheng, Peijin Han, Xuan Hui, Joseph Moore, Scott
247 Robertson, Charles Mayo, Ranh Voong, Harry Quon. Practical data collection and extraction for
248 Big Data applications in radiotherapy. Medical Physics. (Accepted for PBDW2017 Issue): 2018.
249 (In Press)
250
- 251 11) Bhadrasain Vikram. Perspectives on potential research benefits from big data efforts in
252 Radiation Oncology. Medical Physics. (Accepted for PBDW2017 Issue): 2018. (In Press)
253

254 12) Wei Zou, Huaizhi Geng, Boon-keng K Teo, Jarod Finlay, Ying Xiao. NCTN clinical trial
255 standardization for radiotherapy through IROC and CIRO. Medical Physics. (Accepted for
256 PBDW2017 Issue): 2018. (In Press)

Author Manuscript