



Bringing Modern Machine Learning into Clinical Practice Through the Use of Intuitive Visualization and Human–Computer Interaction

Richard Osuala¹ · Jieyi Li¹ · Ognjen Arandjelovic¹

Received: 31 July 2018 / Accepted: 23 January 2019 / Published online: 19 February 2019
© The Author(s) 2019

Abstract

The increasing trend of systematic collection of medical data (diagnoses, hospital admission emergencies, blood test results, scans, etc) by healthcare providers offers an unprecedented opportunity for the application of modern data mining, pattern recognition, and machine learning algorithms. The ultimate aim is invariably that of improving outcomes, be it directly or indirectly. Notwithstanding the successes of recent research efforts in this realm, a major obstacle of making the developed models usable by medical professionals (rather than computer scientists or statisticians) remains largely unaddressed. Yet, a mounting amount of evidence shows that the ability to understand and easily use novel technologies is a major factor governing how widely adopted by the target users (doctors, nurses, and patients, amongst others) they are likely to be. In this work we address this technical gap. In particular, we describe a portable, web-based interface that allows healthcare professionals to interact with recently developed machine learning and data driven prognostic algorithms. Our application interfaces a statistical disease progression model and displays its predictions in an intuitive and readily understandable manner. Different types of geometric primitives and their visual properties (such as size or colour) are used to represent abstract quantities such as probability density functions, the rate of change of relative probabilities, and a series of other relevant statistics which the healthcare professional can use to explore patients' risk factors or provide personalized, evidence and data driven incentivization to the patient.

Keywords Health care · Data · Visualization · Medicine · Patient · Interaction

Introduction

Electronic medical records (EMRs)—also referred to digital medical records, or electronic health records—nowadays a routinely collected data resource in hospitals in economically developed countries, offer an exciting opportunity for machine learning-based knowledge discovery which could significantly affect healthcare delivery, its quality, and therefore intervention outcomes [1–5]. Some of the most prominent problems addressed by the existing literature include the discovery of risk factors, the modelling of disease progression patterns, and the development of patient specific prognostics [6–9]. However, a

major challenge posed by the need to interface these technological advancements with medical personnel and patients themselves, has attracted much less research attention [10–14]. Yet, some of the very premises of the work on person specific prognosis include the incentivization of patients [15]. Moreover, the ability to interact with technology in an intuitive manner is a major aspect governing its adoptability in actual healthcare practice [16, 17].

The visualization tools we introduce in this work are built around a recently proposed disease progression model which has demonstrated highly promising results on real-world data [8, 15]. This model, and indeed all models likely to be successful on the task of comorbidity modelling and prediction, is highly technical and in that sense not readily accessible to medical practitioners or patients. A large volume of previous work has shown that this can be a major obstacle in the adoption of technology in the clinical context [10, 16]. Thus, the contribution of this

✉ Ognjen Arandjelovic
ognjen.arandjelovic@gmail.com
http://oa7.host.cs.st-andrews.ac.uk/

¹ School of Computer Science, University of St Andrews,
St Andrews KY16 9SX, UK

paper is a novel framework which makes a major step towards bridging this gap of outstanding practical significance.

Under the Hood: The Underlying Prediction Model

For completeness herein we present a summary of the key ideas of the adopted method. For in-depth technical details, and the related discussion and results, the reader is referred to the original publications [8, 18–20].

The history vector-based sequential prediction model we adopt from the work of Arandjelović [15] treats a patient's medical record as comprising a sequence of hospital admissions $a_1, \dots, a_i, \dots, a_n$ which form a hospital admission history H :

$$H = a_1 \rightarrow a_2 \rightarrow a_3 \rightarrow a_n \quad (1)$$

Each a_i is a discrete event coded using one of a number of standard disease coding schemas, e.g. [21] or one of a number of mostly related alternatives [22]. The most likely follow-up admission a_{n+1}^* is calculated by likelihood maximization from the current history:

$$a_{n+1}^* = \arg \max_a p(H \rightarrow a) \quad (2)$$

The method proposed by Arandjelović represents a history as a fixed length binary history vector $v = v(H)$ over the most common disease diagnoses, where 1 denotes the presence of a specific diagnosis in the history, and 0 absence thereof. The transition probabilities between different history vectors $p(v(H_1) \rightarrow v(H_2))$ are learnt from a training data corpus.

The original model described in [8, 15] facilitates sequential prediction only. In other words, it predicts the next diagnosis for a patient (or, equivalently, provides a probability ranked list of diagnoses) without any associated temporal information, i.e. it is not able to predict the timing of this diagnosis. Herein the original model is further endowed with a temporal predictive ability. This is achieved by learning the cumulative distribution function (cdf) of a transition from one history vector to another. Considering that a appropriate probability density function (pdf) associated with transitions is the log-normal distribution:

$$p_t(t) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{\ln(t) - \tau}{\sigma\sqrt{2}} \right) \right] \quad (3)$$

the corresponding cdf is:

$$P_t(x) = \frac{1}{t\sigma\sqrt{2\pi}} e^{-\frac{(\ln t - \tau)^2}{2\sigma^2}} \quad (4)$$

where t is the temporal distance of the transition measured from the present, and τ and σ the parameters of the distribution, frequently referred to as the 'location' and 'scale' parameters. The two parameters are also readily learnt from the training data corpus using standard maximum likelihood estimation.

Sequential, Non-temporal Visualization

As explained in the previous section, the cornerstone representation in the model which we build our visualization around comprises vectors with binary entries. This conceptually simple representation allowed us to come up with an elegant design which immediately draws the user's attention to salient features in a patient's medical history. The main window of our application is shown in Fig. 1. To start, consider the bottommost row of filled circles. Each circle corresponds to a diagnosis included in the predictive model, as indicated by the corresponding diagnostic code underneath. Notice that the only aspect in which the appearance of a circle can vary is its colour. In particular, we denote diagnoses present in a patient's history using dark blue and those which are not present using light blue.

Next, observe that there are multiple histories displayed concurrently. The bottommost history, labelled 'Initial History', corresponds to the history from which the space of possible diagnostic trajectories is explored. In clinical practice this initial history will usually be the diagnostic record of a patient at admission. Thereafter exploration proceeds by the user selecting a specific diagnosis (by clicking the corresponding circle). This action changes the history denoted 'Current History' which corresponds to the current state in the exploratory process and is guided by information in the topmost row. Unlike the three other rows which display the same type of information, namely diagnostic histories, the circles in this row also vary in their size and colour which encode the probability of a specific diagnosis given the current diagnostic history, estimated using the model detailed in the previous section. Thus, the user is informed in the exploratory process and can pursue possible diagnostic futures which are more likely.

Selective Emphasis

Recall that the original work which introduced the adopted prediction model based on the diagnostic history vector on the 30 most common diagnoses. This is a sufficiently low number to allow for the visualization described in Sect. 3 to appear uncluttered on most devices. However, subsequent work has demonstrated that the model is successful even with the inclusion of a much greater number of

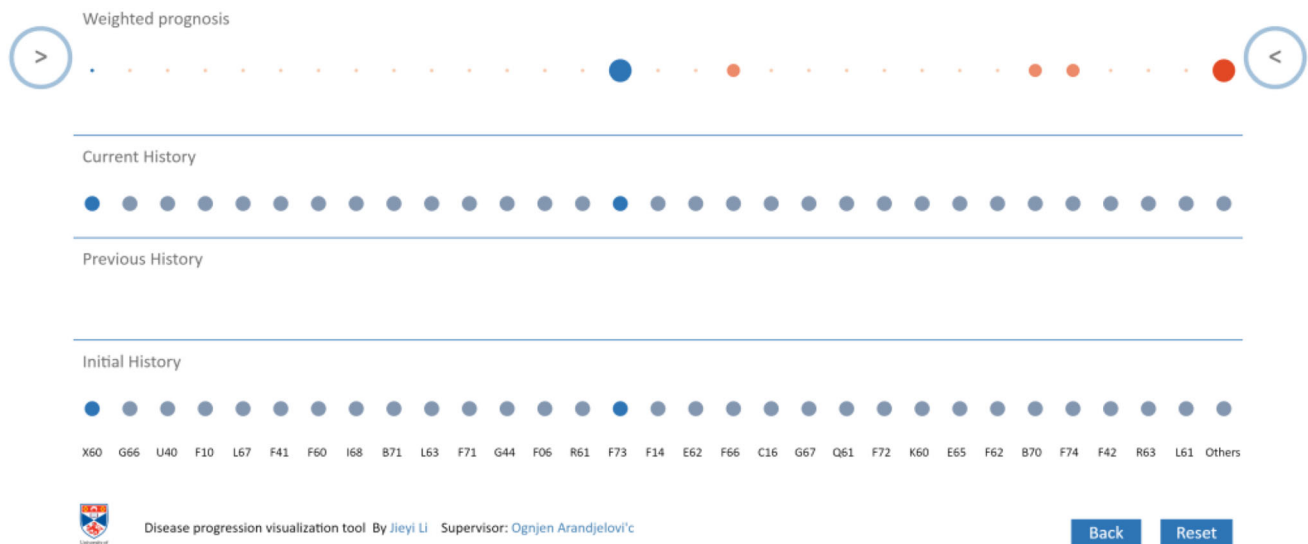


Fig. 1 The main window of our first visualization tool which uses only a sequential prediction model rather than one endowed with probabilistic temporal information (see next section)

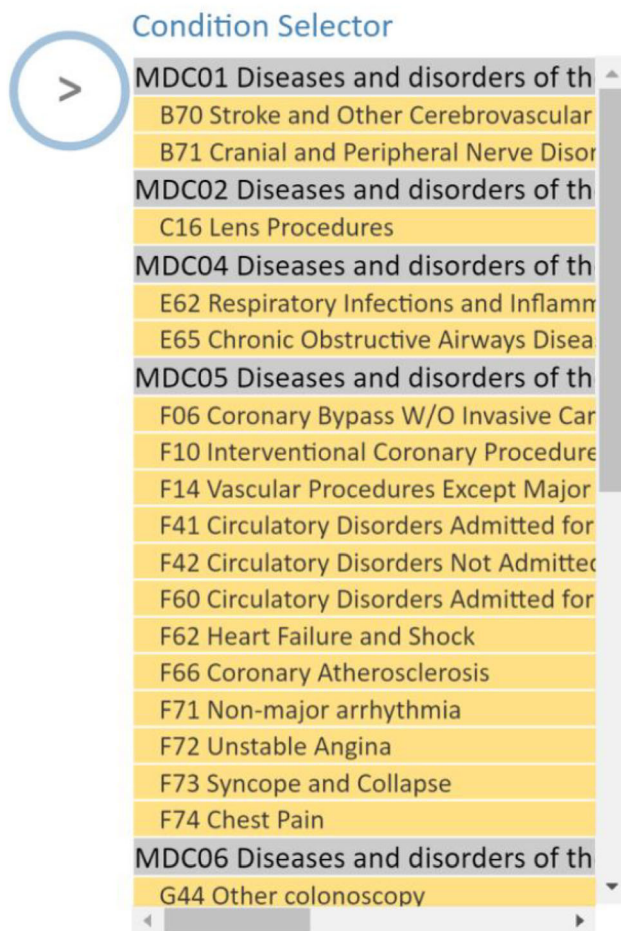


Fig. 2 Selector of diagnoses of interest. Diagnoses can be selected or deselected for display purposes, or grouped according to the hierarchy of the used diagnostic coding schema

diagnoses which can be of clinical interest [19]. Attempting to visualize these in the same manner clearly poses problems with clutter and the ease with which diagnoses of interest can be observed. To overcome this obstacle, we came up with two solutions. Firstly, we allow the user to select or deselect specific diagnoses from being visualized (Fig. 2). Deselected diagnoses are still included in the predictive model, but their states are not displayed in the main window. Secondly, we make use of the hierarchical nature of diagnostic coding. In particular, our application supports several common coding schemes, including ICD-10 and AR-DRG, and thus allows for diagnoses to be grouped according to the subtree in the hierarchy. In other words, rather than displaying related diagnoses separately, the user can choose to unify these and visualize merely that any of the diagnoses of a specific group are present. As before, if more granularity is desired at any point, the option can be changed and individual diagnoses displayed, given that it is only the visualization which is altered and not the underlying predictive model.

Additional Prognostics

The exploration of diagnostic futures described thus far is local in the sense that the user can see the predictions of short-term risks and using this information make incremental moves through the tree of different possibilities. However, considering that the original work has demonstrated good performance on the task of long-term prediction, we also sought ways of visualizing this aspect of the model too. This information is useful in that it can be time saving, more incentivizing to patients, and direction providing in incremental exploration. Hence, we provide

the option to display an additional type of prediction. In particular, we sample ultimate diagnostic histories reached from the current history and display the probability of each diagnosis according to the proportion of ultimate histories in which it appears. As before, probabilities are encoded using size and colour, as shown in Fig. 3.

Inclusion of Temporal Information

The tool described in the previous section allows for the visualization of sequential information only without any associated temporal understanding. Yet, in the present context time is critical – it is necessary for stratifying patients into low and high risk categories, and for allocating resources. However, the incorporation of temporal information in an easily understandable manner is challenging. In addition to the most obvious information which is ‘time until event’ (or rather, the probability density function corresponding to it), the rate of change in instantaneous risk is of importance, and different temporal characteristics of comorbidities can effect a sequencing change over time, which are factors which all add to the complexity and multinationality of information which needs to be displayed. By consulting with a number of relevant healthcare professionals (clinicians, doctors, and nurses) and by adopting an iterative design-test-reassess design process, we found that different users found different manners of information presentation most intuitive and easiest to understand. Consequently we developed a combination of different visualization options which can be readily switched between by the user.

Blob-Based Visualization

The first circle-based visualization approach resembles a so-called blob chart [23], with equidistant blobs which represent different disease diagnoses being distributed horizontally, as illustrated in Fig. 4. The corresponded

diagnoses are labelled using their codes under the adopted coding system (e.g. WHO’s diagnosis-related groups [21], or the Australian refined diagnosis-related groups [22]). As noted earlier, these are standard codes, used widely and understood by healthcare professionals, and allow for the diagnoses to be shown in a succinct, clutter free manner. Additional information and a more detailed description of a diagnosis can be obtained by clicking any visualization element associated with the diagnosis (its label or the corresponding blob, in this visualization).

The size of a particular blob encodes the value of the cumulative density function corresponding to the occurrence of the respective diagnosis by the specific time in future. This time is specified by the user and allows the user to gain an understanding of the highest risks for the patient within this period. Larger and thus more prominent blobs (and hence the corresponding diagnoses) draw the user’s attention to the most probable complications while at the same time providing a simple way of judging relative risks too—several large blobs immediately suggest a cluster of comorbidities, whereas single dominant blob highlights a specific primary diagnosis of interest.

Moreover, we encode the rate of the cdf change by a blob’s colour, using the standard heat map. In this manner, in addition to the instantaneous value of the cdf, we communicate to the user the possibly uneven changes in the probabilities of different diagnoses over time. By including this information in our visualization, a clinician can be alerted of a high risk increase in the near future (relative to the currently selected date of interest). The blob chart visualization is set as the default visualization, whereas the other two visualization options (described next), if selected, are displayed in modal windows.

Bar Chart-Based Visualization

The second visualization option uses the well-known bar chart encoding with the height of bars representing the corresponding value of the cumulative distribution function

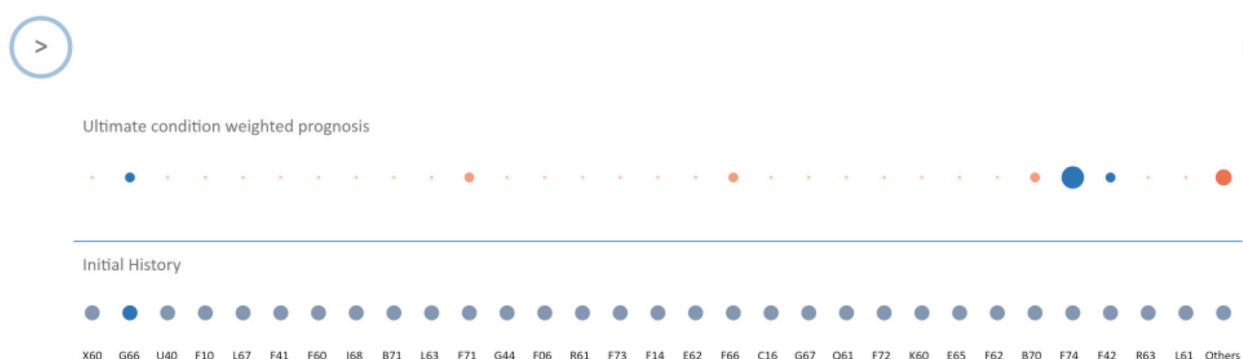


Fig. 3 Ultimate condition weighted prognosis shows the proportion of ultimate histories in which different diagnoses appears

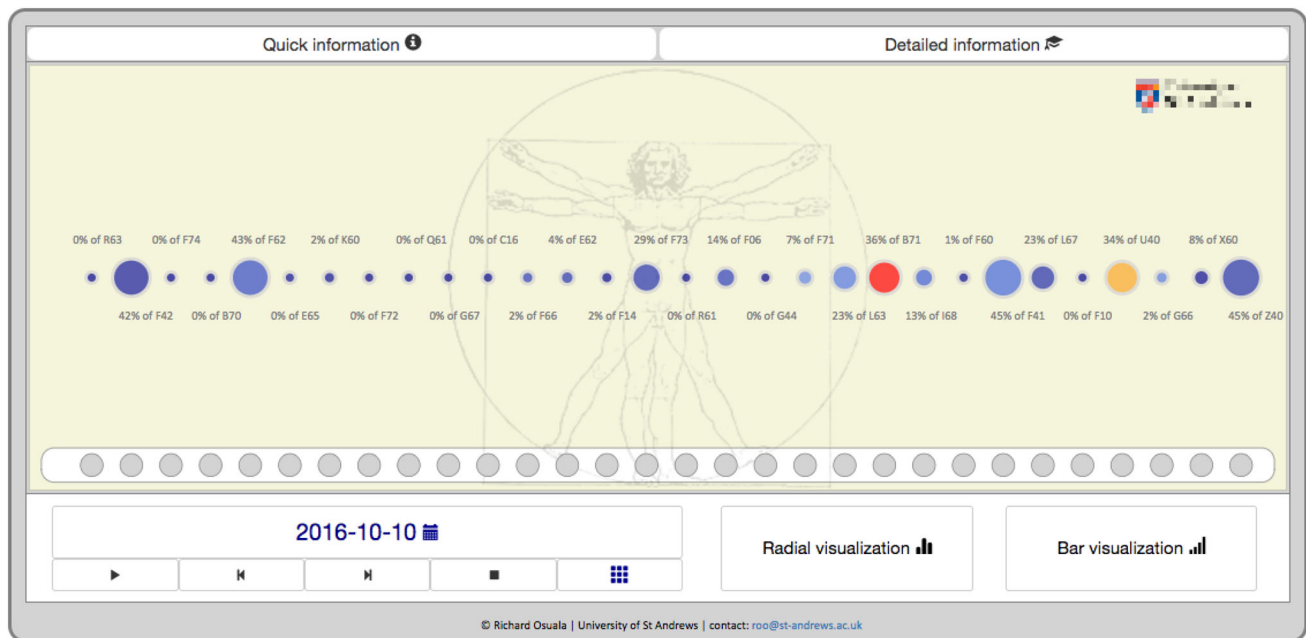


Fig. 4 Blob chart visualization is the default visualization in our application. Each blob represents a disease diagnosis with the blob size encoding the value of the corresponding probability density

function at the selected instance in time. In the example shown in this figure all history vector entries are set to 0, indicating the absence of any diagnoses in the patient's medical history

at the specific date of interest, as illustrated in Fig. 5. As before, bar colours communicate the rate of change of the cumulative distribution function across time and the smaller rectangles underneath the bars represent the presence (or lack thereof) of different diagnoses in the current

history vector. Fundamentally this visualization conveys the same information as the other two alternatives, namely the blob-based and radial chart-based visualizations shown, respectively, in Figs. 4 and 6, but its different way of encoding this information was found to be preferred by

Fig. 5 Bar chart visualization was found to be preferred by some users. Fundamentally it conveys the same information as the other two alternatives, shown in Figs. 4 and 6 but differently encoded

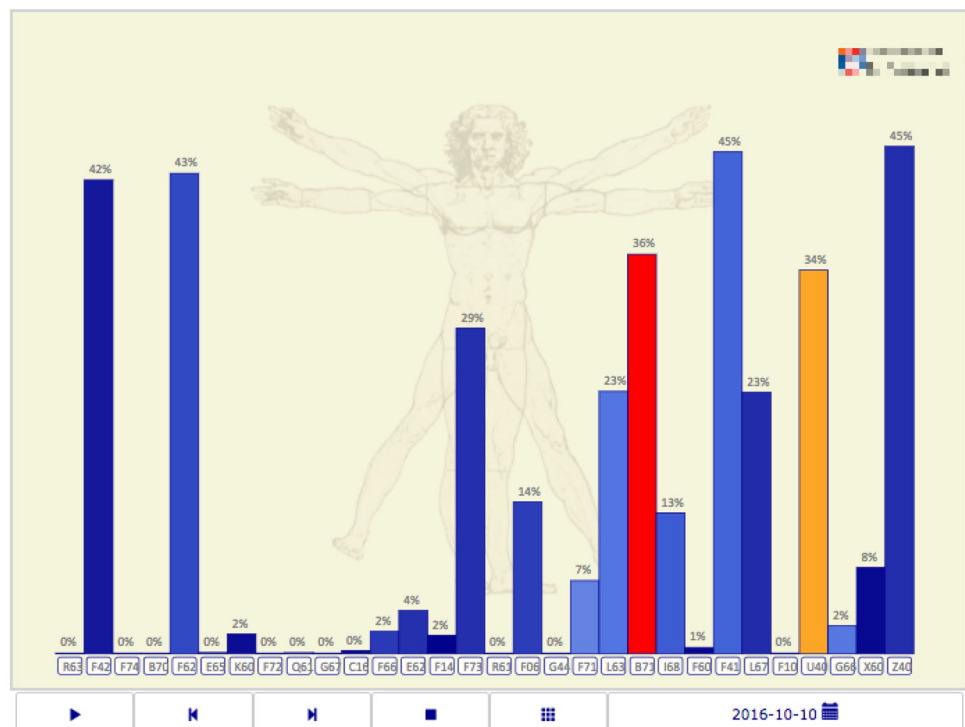
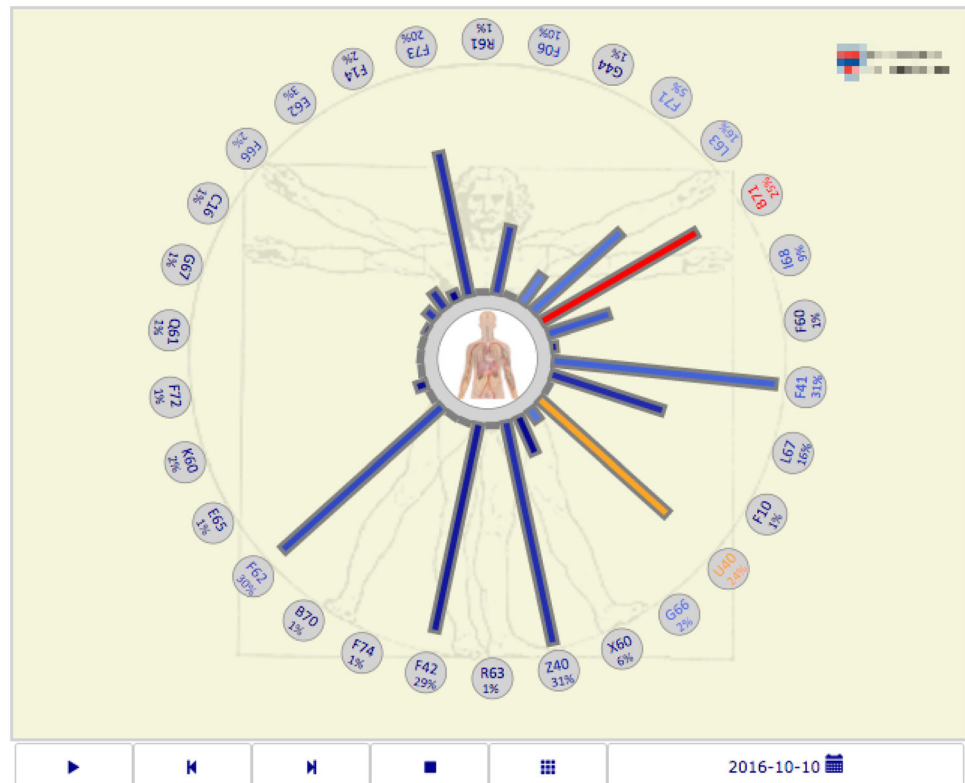


Fig. 6 Radial chart visualization was found to be preferred by some users. Fundamentally it conveys the same information as the other two alternatives, shown in Figs. 4 and 5 but differently encoded



some users. Hence we found it to be a useful alternative to include.

Radial Chart-Based Visualization

The third and final visualization we developed resembles a radial chart with rectangles spreading out radially, as shown in Fig. 6. During the course of our interviews with the target users, we found that some of them preferred this layout to the two described previously due to its symmetry—the lack of symmetry in the first two visualizations suggested to some users some differentiation between different diagnoses which is neither intended nor present in the underlying method or its output. As with the previous visualization, the radial length and colour of rectangles are used to represent the value of the corresponding cumulative distribution functions and their rates of change. We found that with this visualization the users also clearly associated the circles that represent the binary history vector entries (i.e. the presence of specific diagnoses in a patient's EMR history) with the corresponding diagnoses.

Interactive Features

In all of the visualizations, by selecting any cdf encoding element the user can open a window which displays further detailed information and allows for the status of the

diagnosis to be changed, as illustrated in Fig. 7. A fast way of flipping the status of a diagnosis (present or not present) is also provided—a user can simply click on the green (add) or the red (remove) buttons. This effects a history vector transition which in turn triggers a change in the corresponding visual representations (e.g. blob size and colour). This feature allows the healthcare practitioner to explore how different potential diagnoses (e.g. those that the patient may be at the greatest risk of developing) affect the patient's health state further in the future. This can be used as a powerful incentivization tool. For example, the patient can be shown how a specific ailment that he/she is at the risk of developing due to lifestyle choices (e.g. smoking, excessive food intake, etc), would influence other health-related outcomes (e.g. lung cancer, diabetes, hypertension, etc).

Help, Hints, etc.

To facilitate instantaneous help and a ready understanding of different visual elements in our visualizations, when the cursor hovers over any of the relevant geometric entities, all associated information is emphasized. For example, as illustrated in Fig. 7, after hovering over the blob representing the first disease included explicitly in the model [15], a line connecting the blob with the corresponding

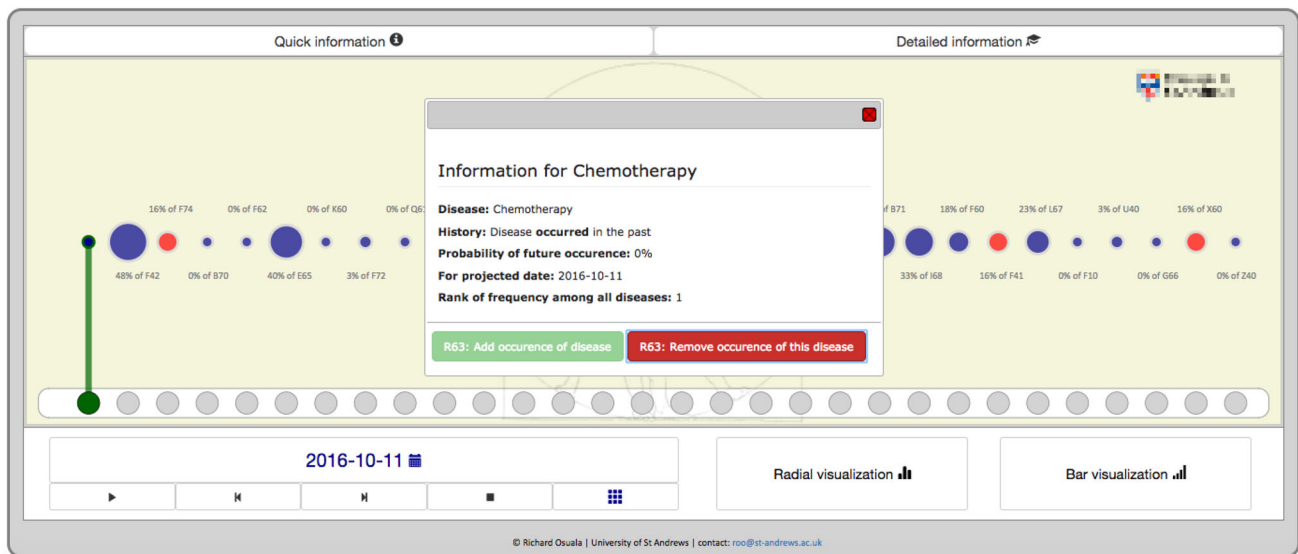


Fig. 7 In this example, a chemotherapy diagnosis is present in the patient's history vector. By selecting the corresponding blob, the user can open a modal window from which the diagnosis status can be

changed. A line connecting the value of the probability density function corresponding to the diagnosis and the date selected (2016-10-11 in this case) is shown so as to emphasize this information

value of the probability density function is shown. Further guidance on the features accessible from the main window of the application is also available, see Fig. 8, as well as a thorough step-by-step tutorial with comprehensive usage information, see Fig. 9. Furthermore, in addition to the disease code, a full description of the diagnosis is shown both above the cursor and at the bottom of the window. To avoid so-called change blindness, further animations emphasize transitions between history vectors or changes

to the date of interest. This feedback uses highlighting and increased contrast against the beige coloured background of the current visualization after 0.5 s.

To minimize the chance of human error and assist the user in interaction, input checking and succinct, timely, and informative notification messages inform users of their interactions with our application. Notifications appear for example in case of validation errors (e.g. date needs to be in the future) or if the date or the history vector is changed

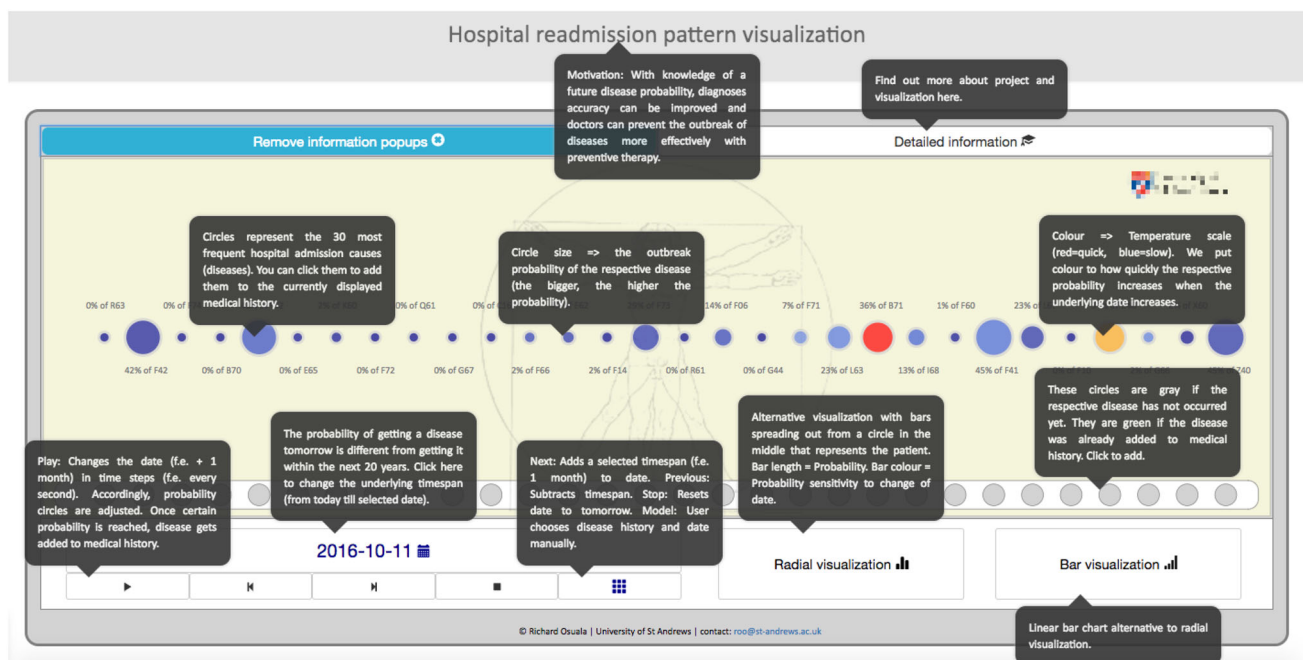


Fig. 8 Succinct and easily understood explanations is readily displayed for all features

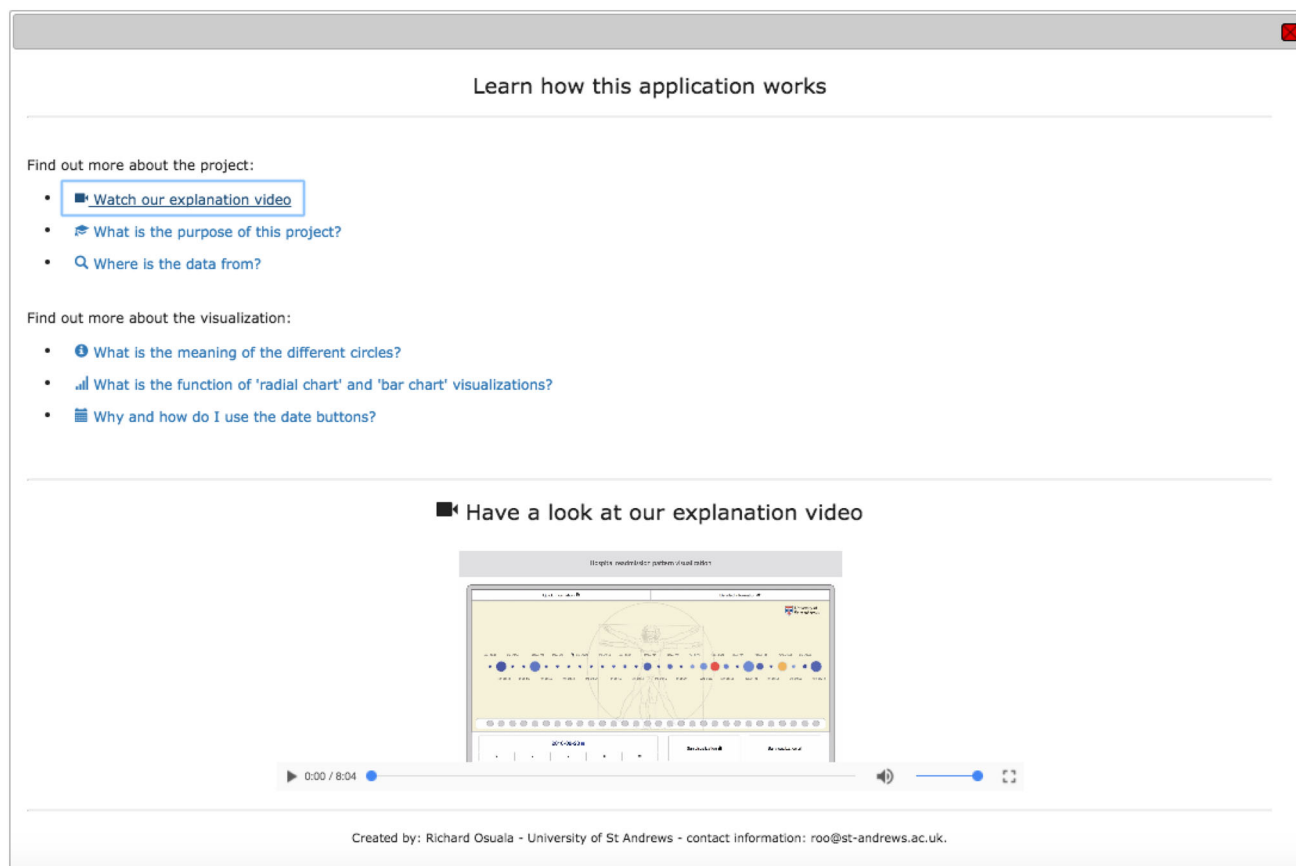


Fig. 9 Comprehensive help and guide for in-depth detail on the use of the application and its features

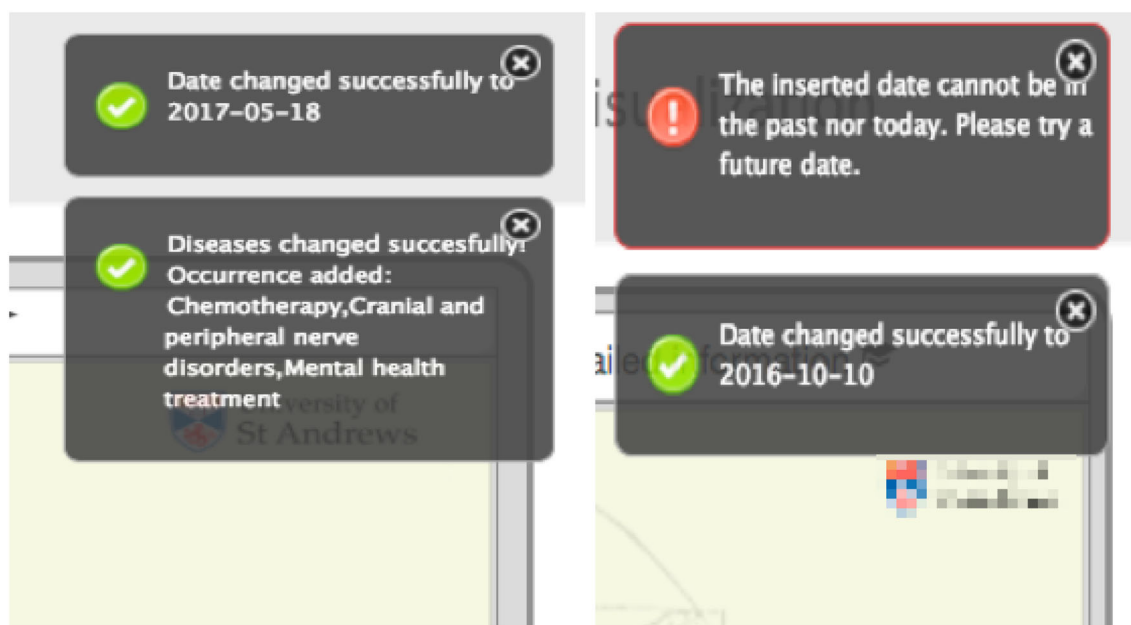
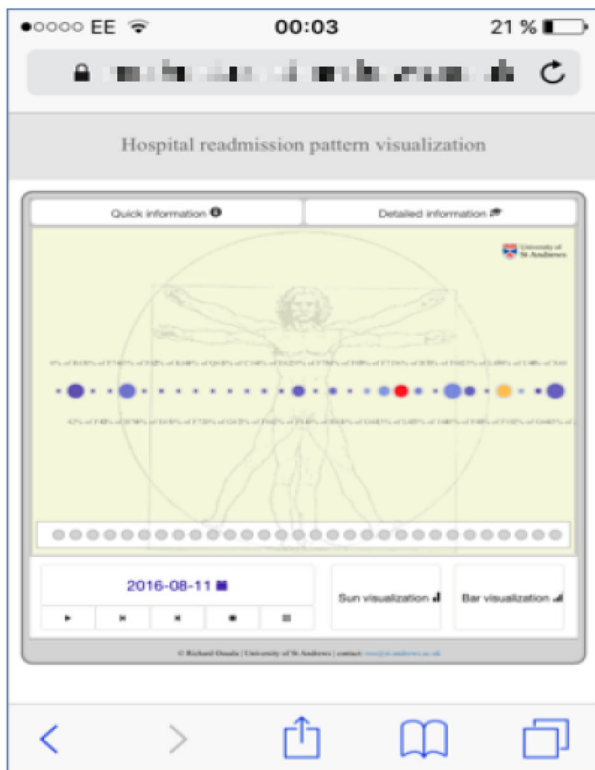


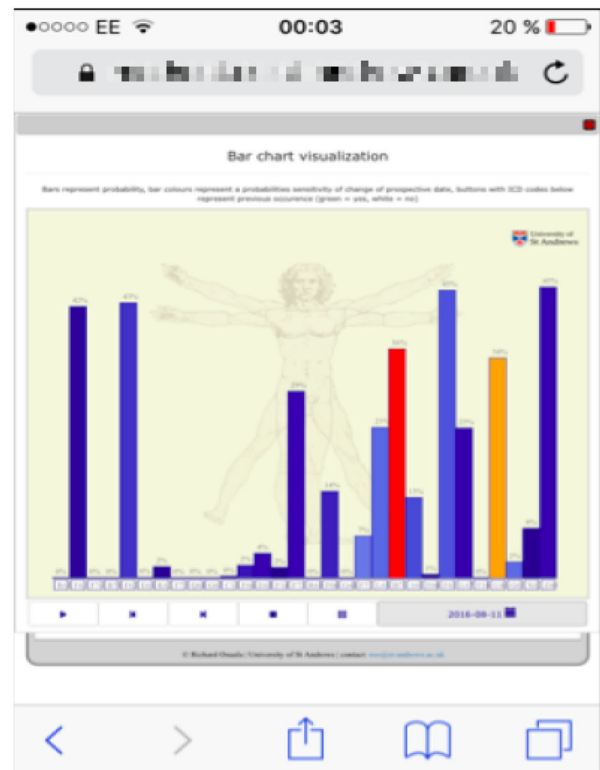
Fig. 10 Examples of notification messages displayed as feedback

due to an added diagnosis, as illustrated with a few examples in Fig. 10. The selected date of interest is visible

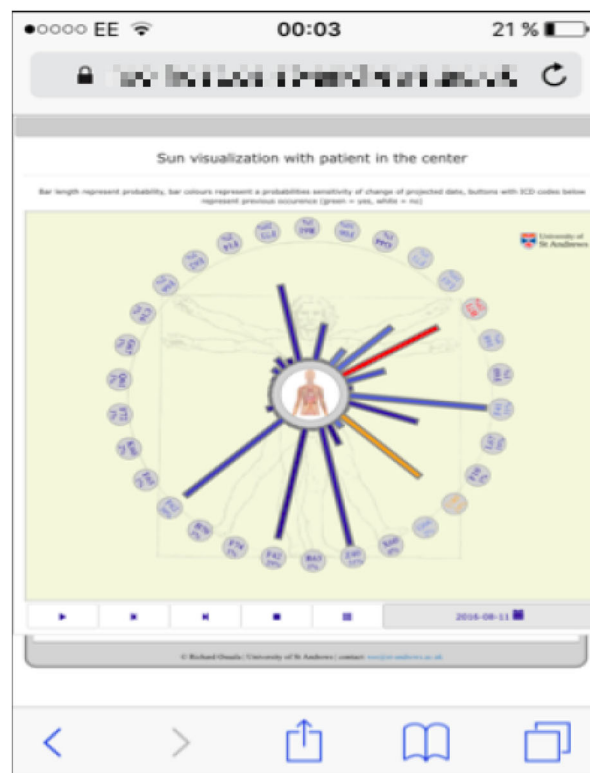
inside an interactive button below the visualization that can be clicked to make adjustments. After clicking on the



(a)



(b)



(c)

Fig. 11 Three visualizations as seen on a 3.5-in. screen in Mobile Safari 9 in an IOS 9.3.1 environment on a cell (mobile) phone

button, a modal window is opened which allows the user to change the date using the familiar calendar view.

Automatic Time Lapse and Long-Term Outcome Simulations

Our application also provides further interactive features, activated using buttons placed below the main visualization space. These buttons resemble the widely known and hence intuitively understandable functions of a media player, such as ‘play’, ‘pause’, ‘stop’. These buttons **provide effortless navigation through time via simulations of possible temporal trajectories through the space of possible diagnoses**. Temporal transitions predicted by the adopted model are accompanied by the automatic visualization of the corresponding disease progression. The forward and backward buttons allow for manual time jumps. Such time jumps change the date of interest and update the visualization accordingly. The duration of such time jumps (e.g. days, months or years) can be specified in the ‘date selection’ modal window.

Clicking the play button opens a modal window where users can also choose the length of time jumps, the real time between predicted transitions (e.g. every 2 s), and whether diagnoses should be added automatically upon exceeding a certain probability of occurrence (i.e. the corresponding cdf value). In the latter case, the play function can add future diagnoses deterministically by using maximum likelihood prediction or non-deterministically by pdf weighted random sampling (ensuring that more likely diagnostic paths are simulated with the correspondingly higher frequency). Once the play function is activated in the modal window, our application repeatedly makes forward temporal jumps (as explained earlier, their duration can be set by the user). If deterministic prediction is selected, diagnoses are added to the visualized medical history if the corresponding cdf exceeds a probability threshold which too can be specified by the user. Random sampling adds diagnoses using cdf-based weighting, thus allowing the clinician to explore multiple future disease progression patterns with repeated activation of the function. When the play function is running, for clarity the ‘play’ button disappears, and is replaced by the ‘pause’ button. The click event of the pause button puts the play function on hold to enable users to explore the currently displayed simulated healthcare record in detail. Clicking the stop button terminates the play function and resets the date of interest to its default value (the present date).

Note on Implementation

Our visualization was implemented as a web application using the D3 Javascript library d3.js thereby offering high

portability across different devices and operating system environments; see Fig. 11 for an example. Additional advantages offered by its web-based implementation include the simplicity of deployment, as no installation or configuration is needed, and an immediate sense of familiarity for non-technical users.

The d3.js-based circles and rectangles used to visualize blobs and bars are nested in a scalable vector graphics (svg). Their radii and lengths are calculated using d3.js scale functions. Heat map colouring uses chroma.js interpolation between four plain colours and scaling with d3.js to calculate the corresponding mapping between the pdf rate of change values and the computed colour palette. To switch from the default blob chart to another visualization format, JQueryUI-based modal functions append HTML code to the interface.

Summary

In this paper we introduced an intuitive visual interface built around a recently proposed computational model of disease progression, aimed at making the model’s predictions accessible to health professionals in their daily work. A range of interactive features allows the user to explore patient specific risk across time. To the best of the authors’ knowledge, this is the first attempt at bridging the gap between increasingly complex machine learning-based algorithms and the realm of healthcare practice. **We trust that our contribution will facilitate increased adoption of technology in healthcare delivery, empowering both the medical community and patients in understanding risk and how to address it. Moreover, we hope that our work will inspire future research in this realm.**

Compliance with Ethical Standards

Conflict of interest On behalf of all authors, the corresponding author states that there is no conflict of interest.

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

References

1. Zhou S-M, Fernandez-Gutierrez F, Kennedy J, Cooksey R, Atkinson M, Denaxas S, Siebert S, Dixon WG, O’Neill TW, Choy E, Sudlow C, Brophy S (2016) Defining disease phenotypes

- in primary care electronic health records by a machine learning approach: a case study in identifying rheumatoid arthritis. *PLoS ONE* 11(5):e0154515
2. Lau EC, Mowat FS, Kelsh MA, Legg JC, Engel-Nitz NM, Watson HN et al (2011) Use of electronic medical records (EMR) for oncology outcomes research: assessing the comparability of EMR information to patient registry and health claims data. *Clin Epidemiol* 3:259–272
 3. Nadkarni PM (2010) Drug safety surveillance using de-identified EMR and claims data: issues and challenges. *J Am Med Inform Assoc* 17(6):671–674
 4. Li J, Arandjelovic O (2017) Glycaemic index prediction: a pilot study of data linkage challenges and the application of machine learning. In: *Proceedings IEEE international conference on biomedical and health informatics*, pp 357–360
 5. Yue X, Dimitriou N, Arandjelovic O (2019) Colorectal cancer outcome prediction from H&E whole slide images using machine learning and automatically inferred phenotype profiles. In: *Proceedings international conference on bioinformatics and computational biology*
 6. Bartolomeo N, Trerotoli P, Moretti A, Serio G (2008) A Markov model to evaluate hospital readmission. *BMC Med Res Methodol* 8(1):23
 7. Duffy ND, Yau JFS (1995) Estimation of mean sojourn time in breast cancer screening using a Markov chain model of both entry to and exit from the preclinical detectable phase. *Stat Med* 14(14):1531–1543
 8. Vasiljeva I, Arandjelović O (2016) Towards sophisticated learning from EHRs: increasing prediction specificity and accuracy using clinically meaningful risk criteria. In: *Proceedings of international conference of the IEEE engineering in medicine and biology society*, pp 2452–2455
 9. Neofytos D, Arandjelovic O, Harrison D, Caie PD (2018) A principled machine learning framework improves accuracy of stage II colorectal cancer prognosis. *NPJ Digit Med* 1(1):52
 10. Le T, Reeder B, Thompson H, Demiris G (2013) Health providers' perceptions of novel approaches to visualizing integrated health information. *Methods Inf Med* 52(3):250–258
 11. Nielsen CB (2016) Visualization: a mind-machine interface for discovery. *Trends Genet* 32(2):73–75
 12. Barracliff L, Arandjelović O, Humphris G (2017) Can machine learning predict healthcare professionals' responses to patient emotions? In: *Proceedings of international conference on bioinformatics and computational biology*, pp 101–106
 13. Osuala R, Arandjelovic O (2017) Visualization of patient specific disease risk. In: *Proceedings IEEE international conference on biomedical and health informatics*, pp 241–244
 14. Li J, Arandjelovic O (2017) Intuitive and interpretable visual communication of a complex statistical model of disease progression and risk. In: *Proceedings international conference of the IEEE engineering in medicine and biology society*, pp 4199–4202
 15. Arandjelović O (2015) Prediction of health outcomes using big (health) data. In: *Proceedings of international conference of the IEEE engineering in medicine and biology society*, pp 2543–2546
 16. Bautista JR, Lin TT (2016) Sociotechnical analysis of nurses' use of personal mobile phones at work. *Int J Med Inform* 95:71–80
 17. Thuemmler C, Lim AK, Holanec I, Fricker S (2015) A methodology to assess social technological alignment in the health domain. *IRBM* 37(4):232–239
 18. Arandjelović O (2015) Discovering hospital admission patterns using models learnt from electronic hospital records. *Bioinformatics* 31(24):3970–3976
 19. Vasiljeva I, Arandjelović O (2016) Prediction of future hospital admissions—what is the tradeoff between specificity and accuracy? In: *Proceedings of international conference on bioinformatics and computational biology*, pp 3–8
 20. Vasiljeva I, Arandjelović O (2016) Automatic knowledge extraction from EHRs. In: *Proceedings of international joint conference on artificial intelligence workshop on knowledge discovery in healthcare data*
 21. World Health Organization (2004) *International statistical classification of diseases and related health problems*, vol 1. World Health Organization, Geneva
 22. Kobel C, Thuilliez J, Bellanger M, Pfeiffer K-P (2011) DRG systems and similar patient classification systems in Europe. *Diagnosis-Related Groups in Europe: moving towards transparency, efficiency and quality in hospitals*, 1st edn. Open University Press and WHO Regional Office for Europe, Buckingham, pp 37–58
 23. Munzner T (2014) *Visualization analysis and design*. CRC Press, Boca Raton

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.