



Using Artificial Intelligence to Support Healthcare Decisions

A Guide for Society



Lloyd's Register Foundation
Institute for the Public
Understanding of Risk



KPC4IR
Korea Policy Center for the
Fourth Industrial Revolution



SENSE
about SCIENCE
Because evidence matters



This guide was created through a partnership of: Lloyd's Register Foundation Institute for Public Understanding of Risk, a research institute at the National University of Singapore committed to improving lives by transforming risk communication and the public understanding of risk in Asia and internationally; the Korea Policy Center for the Fourth Industrial Revolution, a research institute at KAIST working to understand and shape emerging technologies and governance of the Fourth Industrial Revolution for a better and inclusive digital era; and Sense about Science, an independent charity that promotes the public interest in sound science and evidence.

We are grateful for the input and personal time given to us by the many data scientists, doctors, researchers and members of the public who were involved in the development and testing of the guide.

Why we need this guide

Artificial intelligence (AI) is software that can use large amounts of data to assess and make predictions – things that human ‘computing power’ can’t do at all or can’t do quickly and accurately. It is ‘intelligent’ because it works out patterns in the data and tests them, rather than just identifying what it is instructed to find – for example, finding patterns in genomic data that might predict who gets a disease, where humans don’t yet know what to look for.

In healthcare, AI has made advances in analysing data about how diseases progress. It is also being used to identify molecules that could make new drugs, diagnose medical conditions more precisely, predict how patients will respond to treatment, and improve the planning of resources such as hospital beds.

COVID-19 has sped up the introduction of these new health technologies. For example, the BenevolentAI platform took one weekend to identify a drug that could be used to treat the new disease—conventional drug discovery methods would have taken eight years.¹ But this rapid introduction of technology has come with the trade-off of less time for robust testing.

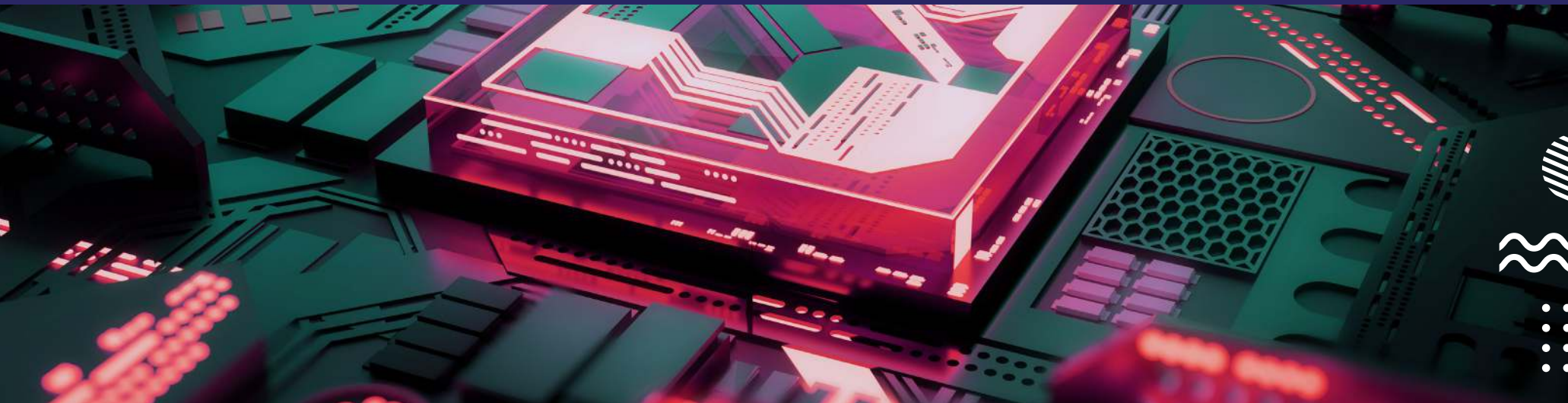
With AI development happening so rapidly, and healthcare providers using AI more and more, it’s vital that more people know the important questions to ask about how reliable different applications are – the quality of the data they are based on, and whether we can depend on them to be right.

It is important for society to ask these questions to make sure AI is used responsibly. This kind of accountability makes a difference: patients asking questions about evidence and outcomes has improved many aspects of healthcare.

Similarly, doctors and patients need to understand how reliable their AI-based information is when life-changing decisions are being made.

But what if policymakers, healthcare agencies, journalists, doctors, and patients don’t know the questions to ask about whether a new breakthrough AI application is reliable or suitable for a particular use? What if they pass on flawed information or make bad decisions because they don’t know where to find information about the model the AI is using? Who is accountable if things go wrong?

This guide is not intended to train AI experts or show how interesting AI is, but to help with the important conversations about its use in healthcare. The guide is designed to equip patients, policymakers, journalists, clinicians and decision-makers with the questions for discussing whether a technology is robust enough for its intended use. It aims to transform the conversation about AI from a complex and daunting one to an empowering one – one that can give us confidence in those technologies that do improve medical treatment and avoid harm from those that don’t.



Terms p06

Technical terms aren't needed to ask the right questions. But where they are used, it helps to know that terms like "AI", "algorithm", "reliability", "model" and "generalisability" have specific meanings.

How AI is used in treating patients p08

AI is helping medical professionals in some fields to work more quickly and accurately, but it can't replace the doctor. Good use of AI depends on its suitability for the decision and the expertise of the medical professionals interpreting it.

Reliability matters p10

There is a lot at stake. AI can base its recommendations on false or misleading relationships it finds in the data, leading to bad decisions. It can make biases in healthcare worse if the limits of the data are not clear. We can only know how reliable AI is if its testing and performance are clear. Understanding how to check on this is important for journalists who want to report on new developments responsibly. It helps health authorities to select the applications that genuinely improve patient treatment, and it helps the public to have confidence in the right things.

Questions to ask about AI in healthcare p12

What data is it based on?

To reduce the chance of the AI identifying false or misleading relationships, it's important to know how the data underpinning it was generated.

What assumptions is AI making about patients and disease?

An AI-supported diagnosis or treatment option might not be useful if the results can't be generalised across countries or groups, or if key information is missing.

How much decision weight can we put on it?

AI can only support a clinical decision if we know how well it performs.

A reliable future p22

To make sure we identify genuinely useful innovations, we must ask the right questions now about the reliability of the AI being used for different purposes. The questions in this guide will help society create a benchmark for responsible discussion, that will promote clarity and high standards for the use of AI in healthcare.

Terms



Algorithm

A set of mathematical instructions to find or calculate something. Algorithms can be used by AI to find relationships between things (variables) in data.



Artificial intelligence (AI)

A machine or system that uses data and rules to make assessments or predictions like a human would.



Model

An equation that an AI uses to represent how conclusions can be made from data the AI hasn't seen before. For example, new information about changes in smoking habits can be used in a model to predict the number of cases of lung cancer.



Reliability

How trustworthy an AI is or how consistently an AI produces the result we want (e.g. being better at identifying the patients whose disease will improve with surgery) without producing results we don't want.

It can also mean, technically, the ability of an AI to produce the same result every time.



Big data

A type of data that is large (volume), varies in content and type (variety), and changes quickly (velocity).

In the healthcare context, such data includes many variables (e.g. age, gender, height, weight, average weekly alcohol consumption, smoking habits, chronic conditions, medical treatments, test results and x-rays) and can be in different formats (e.g. sounds, videos, written records, images, charts and graphs).



Generalisability

A measure of whether the conclusion made using a set of data is generally true or not. For example, an AI that is not generalisable can help with a diagnosis of bone conditions for only certain demographic groups but not others.



Variable

A factor or characteristic that might be relevant to answering a question. These could be numbers like age, weight, height, temperature or income. Or they might fall into categories like eye or hair colour, ethnicity, field of work or hobbies.



How AI is used in treating patients

AI is intended to help medical staff work quickly and accurately and to make processes efficient.

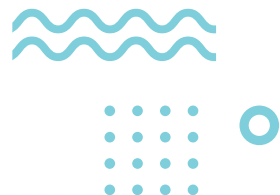
Current AI-based software is limited to performing specific tasks to support a doctor's decision-making. It cannot perform complex tasks such as making clinical decisions, and doctors can consider things that the AI cannot, such as a patient's cultural practices, when making a treatment plan.

At the current pace of technological development, this is likely to be the case for the near future: AI can support but not replace the doctor.

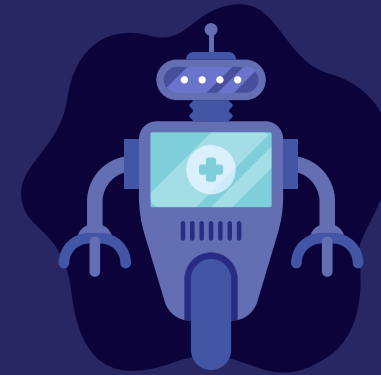
In South Korea, VUNO Med solutions are AI-based diagnostic support systems that can read medical images or analyse biosignals. VUNO's BoneAge assessment software compares bone age with chronological age - for example, an eight-year-old child whose bone age is nine years old is assessed to be growing too fast.²

In Germany, a diagnostic AI has been used to detect potentially cancerous skin lesions. It was tested against an international group of 58 dermatologists and proved better at correctly identifying the nature of more suspicious lesions.³

On the other hand, an eye disease diagnostic developed by Google Health⁴ suffered from a major drawback: the quality of many images taken by nurses was not high enough, so the system rejected more than a fifth of the images and more work had to be done to retake these images. The theoretical accuracy of the diagnostic prediction can only be realised if medical professionals have the confidence and training to use it.



Types of AI in healthcare



Clinical-decision support tools

Medical devices and applications used by clinical practitioners to perform their work. AI is used in diagnostic imaging, predicting treatment outcomes, robotics in surgery and remote monitoring of patients who are using medical devices.

Patient-decision support tools

Medical devices and applications used directly by patients or caregivers. Examples include chatbots or other online tools which help with self-diagnosis, and lifestyle applications such as fitness trackers.

Healthcare administration

Tools used by organisations to improve operations and administration - AI is used in resource allocation, cost reduction (e.g. by reducing test duplications) and automating processes like dispensing medicines.

Therapeutics development

AI used in discovering new drugs and treatments.

Reliability matters

The use of AI to help with diagnosis, predict the outcome of treatment or prioritise resources is potentially life-changing.

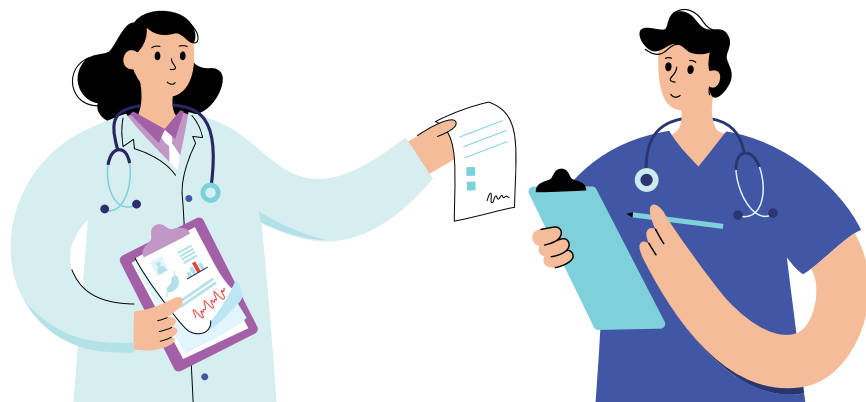
There is some suspicion about AI among the public and healthcare practitioners. Its inner workings are difficult to see, which makes it difficult to question or contest its conclusions, and there are fears about how it uses personal data.

Privacy issues are often raised, but reliability issues have been neglected, perhaps because it is difficult to know how to question them. While it is important for people to have confidence that their data is secure, it is just as important to know whether data is being used well. It's unlikely that any of us would accept a technology based on a study with a 10-person sample size on the basis that the 10 participants' data was kept safe.

Guarantees about privacy are not enough for a technology to be useful, so key questions about the quality of data and reliability of AI need to be asked.

Poor-quality data (or poorly understood data) affects the accuracy of AI. Biases in AI arise from missing or excluded data, existing bias in the training data or errors in the algorithm. Like other data analytics, using data for a purpose it wasn't collected for can introduce false or misleading relationships. And we can't be sure how reliably the AI performs if the model hasn't been rigorously tested in the real world.

The use of AI to help with diagnosis, predict the outcome of treatment or prioritise resources is potentially life-changing.



So, scrutinising quality and reliability means checking that:



The source of the data is known



The data has been collected or selected for the purpose it's being used for



Limitations and assumptions for that purpose have been clearly stated



Biases have been addressed



It has been properly tested in the real world

How do we know that someone has done these checks? There are questions that everyone can ask – whether a journalist, policymaker, clinician, patient or relative – to find out. These questions are set out in the next few sections.

QUESTIONS TO ASK ABOUT AI IN HEALTHCARE

What data is it based on?

Data is obtained in different ways.

Experimental data (collected from experiments) is collected to answer a specific question. Researchers usually consider possible biases they will get in the data and what might be missing, and take steps to overcome these issues.

Observational data is recorded as we go about our business, such as withdrawing money from a bank or travelling on public transport, and there is also administrative data that is recorded by institutions, such as speeding fines or the issuing of prescriptions at hospitals. The biases and limitations in such data sources are usually not thought about until the records come to be used as data for analysis.

All these data sources can be useful for developing AI, but it's important to consider how good and relevant they are for a particular purpose, especially if they've not been gathered for that purpose.

For instance: 'What factors cause patients who have recovered from alcohol addiction to relapse?'

Programmers might put together databases containing information (variables) such as age, chronic medical conditions and genetic profiles.

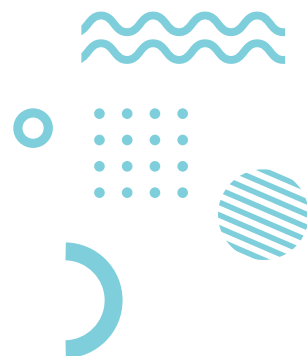
AI would look at these detailed datasets for relationships with relapse data.

If the data came only from medical sources, the AI could miss potential major factors such as unemployment and miss people who do not engage with medical services.

So, the aspects of this question to consider are:





- ✔ **How the data used to train the AI was collected**
- ✔ **Whether the data represents the patients for whom the AI is being used**
- ✔ **Whether the patterns and relationships identified by the AI are accurate**

Not everyone will be able to ask about or assess the details of these aspects, but any doctor, patient or reporter can insist on a clear statement of how these aspects have been addressed. Anyone commissioning the AI for use in health services should be confident that they know the answers.



How was the data used to train the AI collected?

If the data comes from an experiment, it should have been collected to answer a specific research question as part of a well-designed study. Signs of quality include:

-  **A large sample size of participants**
-  **A control group of participants with similar characteristics to compare results against (except for the variable being measured)**
-  **Error estimates**
-  **A discussion of how well the research findings can be extrapolated to real life**

AI systems trained using this type of data have a lower risk of having false or misleading relationships if those quality markers are there.

Observational data analysis involves looking at data that already exists and searching for relationships between variables. There are advantages to this approach, such as being able to study many more variables than an experiment would allow. While it is possible to correctly identify relationships with this type of data, the data source should be clearly stated, and information provided about the AI should include how biases have been considered.

We should also note if the data gathered consists of objective measurements (e.g. vital signs from a device) or subjective self-reported data (e.g. survey responses). Subjective data could have more inaccuracies or biases as people's responses vary for different reasons and responses are self-selecting.

Singapore's Health Promotion Board is collaborating with Apple on an app called LumiHealth. Developed in close collaboration with doctors and public health experts, LumiHealth aims to deliver personal health recommendations based on factors such as age, gender and weight. These recommendations are driven by AI using real-world data from users (obtained with consent) and include reminders to go for regular health checkups. By following the app's recommendations, a user can work towards weekly activity goals and participate in challenges that aim to improve sleep habits and food choices.⁵

How personal and relevant a health app's recommendations are depends on how the data behind it is gathered. LumiHealth uses user data carefully selected for relevance. But some apps don't do that. If an app uses observational data from other users of the app to recommend when a person should visit a doctor, the recommendation is likely to be skewed by the fact that healthier people tend to use such apps.



Does the data represent the patients for whom the AI is being used?

Data might not be useful for training an AI if it doesn't represent the target population. It may be missing information about different ethnicities, sexes, and age groups, and in some cases, this missing information has important implications for health. For example, heart problems show up differently in men and women, or the data may be based on people who can afford to seek treatment and therefore biased to the health of wealthier people.

In Germany, a skin diagnostic AI was trained and validated using images obtained primarily from fair-skinned people in the USA, Australia, and Europe. If the algorithm bases most of its knowledge on how skin lesions appear on fair skin, then there's a risk that lesions on patients with darker skin are more likely to be misdiagnosed.

The absence of data from people with darker skin won't make the diagnostic useless, if it can reliably support diagnoses in some people. But clearly the absence of important data about certain ethnic groups should be known about in countries that have multiracial populations, which is the case in many East and Southeast Asian countries.

Overcoming problems with the representativeness of data is a challenge. Some groups are under-represented in health studies so are under-represented in the data.

Privacy regularly comes up in the public conversation about AI and the use of data. People are concerned that their personal medical data could be used to discriminate against them. For this reason, certain categories of information, such as rare or genetic conditions, require strong anonymisation procedures.

Public concern about privacy influences whether people will share data, and this can affect the accuracy of the AI's recommendation by giving it too small a pool of information to draw reliable conclusions from. By being transparent and demonstrating the steps taken to check that the AI is reliable, researchers and developers can help give people confidence about providing their data.

~~~~~

### Public concern about privacy influences whether people will share data, and this can affect the accuracy of the AI's recommendation.

~~~~~

During the COVID-19 pandemic, Singapore rolled out the TraceTogether mobile app for contact tracing. The idea behind it was the exchange of Bluetooth signals between mobile phones with the installed app. Each phone could detect other participating TraceTogether phones nearby. The app estimated the distance between users and the duration of any time spent less than two metres apart. Encrypted records of these contacts were stored on each user's phone for 21 days. An app user identified as having come into contact with a person who had tested positive for COVID-19 could authorise their TraceTogether data to be accessed by the Ministry of Health (MOH). MOH would then decipher the data and get the mobile numbers of the user's close contacts from the previous 21 days to contact-trace them, ask them to isolate, and test them.

TraceTogether was unable to gain public trust. In June 2020, three months into the app's launch, approximately 30% of the population downloaded the application, falling short of the required adoption rate of 50-70% for contact tracing apps to be effective. Many Singaporeans saw the app as a phone surveillance mechanism. By December 2020, nine months after the app's launch, its adoption rate had barely grown. However, through the government's distribution of an alternative – an external device with the same function – Singapore achieved a 70% adoption.⁶



The experience of Singapore's contact-tracing app shows that real-world limits on data can be hugely underestimated by app developers. People are sometimes just not willing to provide the data that will make even a well-designed application work. Some of these concerns may be alleviated by greater transparency in the technology itself, but in other cases they won't. We have to ask instead whether the application is going to be fed with enough relevant data to continue running reliably.

Finding out if the data is appropriate for its intended use helps to reduce the risk of AI systems spotting false and misleading relationships.

Are the patterns and relationships identified by the AI accurate?

Data is fed into an algorithm, which can analyse the data to find patterns between variables. An AI can learn about these relationships between variables as more data is fed, apply these relationships, and adjust them.

We've seen that the ability to quickly spot patterns in data is a key benefit of using AI in healthcare and that it also presents challenges. It's possible that an AI might start to spot patterns that are not relevant.

Population-level data contains information about lots of variables – for example, people's age, gender, ethnicity, marital status, jobs, postcodes, what car they drive, whether they are registered to vote. This type of data is often called "big data". If an AI searches through enough big data, it will inevitably find patterns and relationships between variables that have nothing to do with each other. This is known as data dredging.

It shows the need to have a specific question to answer when the AI is programmed to look for patterns in a dataset, because it's then less likely to come up with random relationships that affect the validity of the model.

To make sure the relationships are real, anyone commissioning an AI for healthcare should ask if it has been trained using big data and how data scientists have identified the variables most relevant to what the AI is going to be used for. Moreover, even AI trained using big data can be rigorously tested using an independent dataset – as explained in the next section, AI providers should make clear whether this has been done.

We've seen that the ability to quickly spot patterns in data is a key benefit of using AI in healthcare and that it also presents challenges.



QUESTIONS TO ASK ABOUT AI IN HEALTHCARE

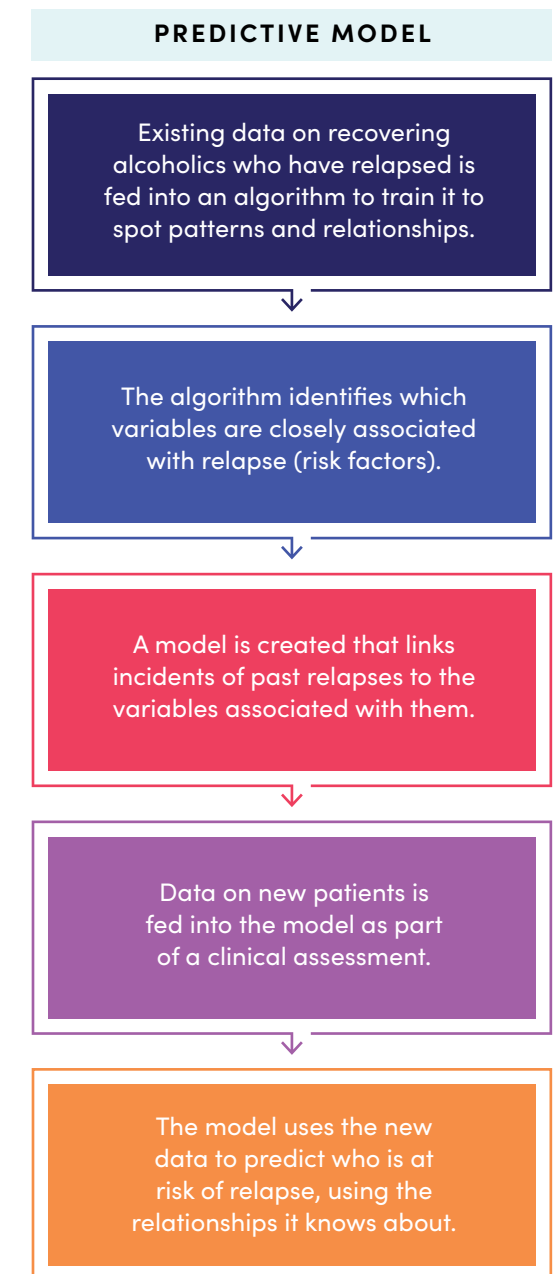
What assumptions is the AI making about patients and disease?

If the algorithm analysing alcohol addiction relapse data finds the relapse rate is higher among low-income people, then it might flag fall in income as a risk factor for relapse. Information about relationships between variables can then be used to create a predictive model – a mathematical equation that uses information about what happened in the past to make a prediction about what could happen in the future (see diagram on the right).

How a model translates to the real world has implications for the reliability, generalisability and fairness of the AI.

The essential aspects are:

- ✔ That the right relationship is captured
- ✔ Whether variables excluded from the model are indeed irrelevant
- ✔ Whether the results are generalisable
- ✔ Whether the AI eliminates human prejudice from decision-making



Is the right relationship captured?

Sometimes, observational data shows up variables that seem to be related to each other (when one goes up, the other goes up or down). Those variables are “correlated”, but that doesn’t mean one “causes” the other.

In 2017, the University of Chicago Academic Hospital System (UCAHS) developed an AI to predict patients’ length of stay.⁷ It was intended to help doctors prioritise patients who were more likely to be eligible for rapid discharge and free up beds faster. The AI’s algorithm found patients’ postcodes to be one of the best predictors for a length of stay. The postcodes associated with a longer length of stay were those in poor neighbourhoods. In effect, the AI recommended

prioritising patients from richer districts. There was a clear correlation between postcode and length of hospital stay, but it is unlikely that a person’s address itself causes them to stay longer in the hospital.

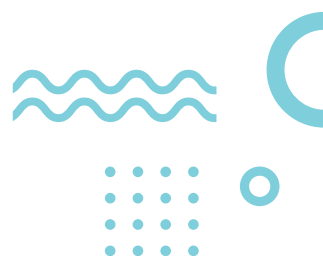
Are the variables excluded from the model actually irrelevant?

The real world contains many millions of variables changing at once. It would be impossible for a model to account for every possible degree of change. Some variables may be readily available and others costly or even impossible to secure. There is only so much computer processing power to draw on and so much time and money to spend. So, data scientists make assumptions and intentionally exclude some variables.

In the Chicago hospital example, the developers needed to consider the missing variables. Perhaps there was a third factor at play, one which caused people to stay longer in a hospital when they are ill. In this case, poverty: poor people live in neighbourhoods where housing is more affordable AND they also tend to have poorer health outcomes and a higher risk of suffering from chronic illnesses. The detrimental

effect of poverty on health is likely to be compounded by a technology that diverts more treatment away from them.

Anyone commissioning an AI product should ask what variables might be missing from the model, why they are missing and how this might affect the outcome. It’s important for developers to understand this themselves and have an open and honest discussion about it with the people they are handing over the technology to.



Are the results generalisable?

AI doesn’t work well when it is required to make a prediction or recommendation on something that differs substantially from its training data (the data used to develop it).

In the VUNO Med-BoneAge example, the AI-based diagnostic supporting solution is built on South Korean population data and validated with multinational data. As the bone growth curve can vary according to race and ethnicity, the accuracy of VUNO’s BoneAge assessment may differ when used with population data from other countries and ethnicities. In this case, fine-tuning or retraining of the algorithm will be necessary to make the output more accurate in Caucasian, Native American and African people.

Variables that could influence the generalisability of an AI application include age distribution, ethnicity, gender, geography and climate. Anyone commissioning an AI product should ask if the results are generalisable, and clinicians should feel confident in the accuracy of the AI’s recommendation for the particular group of patients they are treating.

Does AI eliminate human prejudice from decision-making?

One misconception about AI-supported decision-making is that it is based on cold hard facts without prejudice. But an AI is trained on data from the real world. It sees the world the way it is, not as it could or should be. AI is not inherently good or bad, but it can compound unfairness in healthcare unless the limitations of the data are understood by the developers. Some AI research seeks to address these existing biases through its programming.

If blindly optimising the use of beds was the sole objective, then using patients’ postcodes as a proxy to predict who should be prioritised for treatment wouldn’t be so bad. But the ultimate problem that UCAHS ran into was that poorer people in the USA are disproportionately African American. By diverting treatment away from patients who lived in poorer neighbourhoods, the AI was prioritising white people over black people. It only exacerbated existing racial health inequalities in American society.

Even representative data can embed prejudices, biases and harmful assumptions. In the complex modern world, AI predictions and recommendations can’t be divorced from social realities. Anyone using an AI to aid a clinical decision or any decision in a healthcare setting should consider whether it has the capacity to encode prejudices.

A commentator or patient can ask what assumptions are being made and how we are sure these are fair, even if the AI technically does its job.

This doesn’t mean that any group should be more concerned than another about how AI is used to support their treatment. When the right conversations are had at the right time, everyone involved can be confident in the clinical decision that’s made.

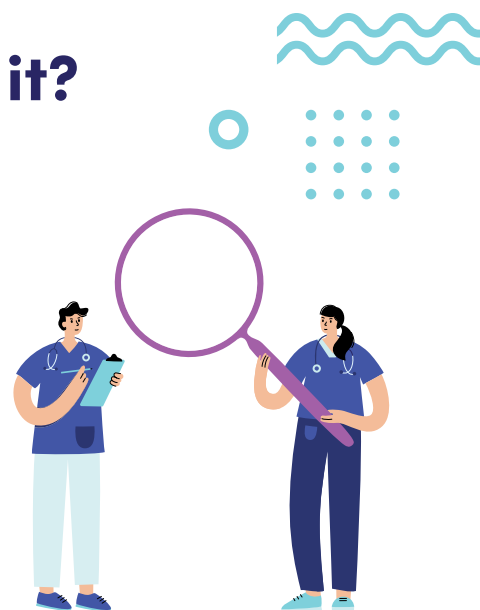
QUESTIONS TO ASK ABOUT AI IN HEALTHCARE

How much decision weight can we put on it?

We've seen that an AI's performance depends on the quality of data it is based on and what assumptions it makes about patients and disease. Taking all this into account makes it more likely that the AI is of good quality, but is it good enough for its intended purpose?

The essential aspects of this question are:

- ✔ **How well the AI really performs**
- ✔ **Whether its reliability has been properly scrutinised**
- ✔ **Whether it makes a useful real-world recommendation**



How well does the AI really perform?

We need to know some basic performance measures that define how good the AI is at predicting things or making recommendations. One measure is accuracy (how often the AI gets its prediction right).

Google Health developed an AI system in Thailand to help identify diabetic retinopathy and speed up the diagnosis process. The process took up to ten weeks while photos of patients' eyes were taken by nurses and dispatched to a specialist for analysis. The AI system could produce results in under ten minutes with 90% accuracy.⁸

But choosing the right way to measure performance is important, and we should be careful not to rely too heavily on theoretical accuracy. With the hypothetical alcohol addiction relapse AI, let's say 10 in 100 recovering alcoholics in this dataset actually relapse after two years. If the AI is 85% accurate at predicting relapses, then it's wrong 15 times out of 100. That means it could miss every relapse and is not much use if it's being used to assess who needs help.

Even if the AI were highly reliable and underpinned by the finest data, a clinician should consider its recommendation in the context of all the other medical evidence they have for a particular diagnosis or treatment option. The doctor makes the final decision.

Has its reliability been properly scrutinised?

As well as accuracy, we should consider the AI's reliability in making predictions. Independent datasets can be used to test how good the AI is at using the relationships it has identified to make a prediction about data it hasn't seen before – its reliability.

This is ideally done by holding back a section of the training data and then seeing how well the AI could identify the thing it's looking for or predict the outcome. Sometimes, an AI that works well on the data used to train it is terrible at making predictions from new data. That could be because the model has not weeded out irrelevant variables or because the model has learnt the training data rather than its underlying relationships. An AI that doesn't make consistent predictions on similar data is unreliable.

Where apps are based on collaboration between public health and the private sector, there is more opportunity to scrutinise reliability. For example, technologies developed in a public-private collaboration are more likely to have undergone clinical trials – set up to see how well it performs against existing practices or human judgement.

Singapore's LumiHealth app was developed via close collaboration between Apple and public health authorities. To be authorised for use in public health, the app needed to meet strict criteria. Close collaboration with public health experts reduced the risk of data not being representative, because the app was not relying on volunteer-contributed datasets.

Does it make a useful real-world recommendation?

One way to determine this is to find out if the AI does any better than a human. It's a good sign if healthcare professionals were involved in the AI's development or deployment. A clinician might look for trials that show whether the AI performs better than, or at least as well as, their trained colleagues.

The German skin diagnostic AI was shown the same images of skin lesions as an international group of 58 dermatologists. It correctly identified the nature of nearly 87% of suspicious lesions compared to 79% for the clinicians. This is one good sign that the AI provided a useful aid to the clinician's decision on treatment.⁹

The AI might also be externally validated, which means tested in the real world. One example would be an AI-based healthcare software company testing its program in a hospital setting to see if it was as accurate in

deployment as it was during testing. The process would be led by experts independent of the AI developers and would show up failures and unintended outcomes.

The process would also identify how the technology would work in practice when subject to human errors in the way it's used: for example, the performance of Google Health's eye disease diagnostic was ultimately hampered by the fact that nurses were not confident in taking high-quality pictures.

We need, finally, to ask what is at stake. A lifestyle app that gives people general advice about diet and exercise perhaps needs only to be roughly reliable. Where the real-world implications of the AI being wrong will be very serious, though, we should expect to see strong evidence of test data, trials and validation.

A reliable future

Using AI to support clinicians in treating patients holds great promise. From rapidly identifying new drug candidates in times of pandemic, to supporting the diagnosis of serious diseases, helping hospitals to manage resources and helping public health agencies to promote healthy lifestyles, AI has demonstrated its value and is here to stay.

But problems arise if the quality of data underpinning the AI is not properly scrutinised and if the AI's reliability hasn't been tested. From misdiagnosing a serious disease to exacerbating racial and economic health inequalities, AI gone wrong can have life-or-death implications. There's confusion and fear out there – fear about robots taking people's jobs, fear about data privacy, fear of who's ultimately responsible if an AI-supported decision turns out to be wrong. Rather than throwing out tools that can help us, we'll be better off if we discuss the right questions now about the standards AIs should meet.

By applying these questions, society can ensure AI developers' solutions to modern healthcare challenges are making good use of the data and knowledge available, with minimal error, across different countries and populations, without deepening inequalities that are already high. These are the AIs that will make useful real-world recommendations that clinicians can have confidence in.

As more people ask the questions in this guide, more people in authority will expect to be asked. In this way, we create a virtuous circle of responsible discussion, and ultimately, higher standards in using AI to guide healthcare decisions.



- 1 BenevolentAI Named as One of Fierce Medtech's Fierce 15 Of 2020 (2020). BenevolentAI. Available at: www.benevolent.com/news/benevolentai-named-as-one-of-fierce-medtechs-fierce-15-of-2020
- 2 VUNO Med@-BoneAge™. Available at: www.vuno.co/en/boneage
- 3 Goyal, M. et al., (2020) Artificial intelligence-based image classification methods for diagnosis of skin cancer: Challenges and opportunities. Available at: www.sciencedirect.com/science/article/pii/S0010482520303966
- 4 Heaven, W.D. (2020) Google's medical AI was super accurate in a lab. Real life was a different story. Available at: www.technologyreview.com/2020/04/27/1000658/google-medical-ai-accurate-lab-real-life-clinic-covid-diabetes-retina-disease/
- 5 LumiHealth™ (2020). Available at: www.lumihealth.sg

- 6 Huang, Z., et al. (2021) Awareness, acceptance, and adoption of the national digital contact tracing tool post COVID-19 lockdown among visitors to a public hospital in Singapore. Available at: www.ncbi.nlm.nih.gov/pmc/articles/PMC7817417/
- 7 Nordling, L. (2019) A fairer way forward for AI in health care. Available at: www.nature.com/articles/d41586-019-02872-2
- 8 Heaven, W.D. (2020) Google's medical AI was super accurate in a lab. Real life was a different story. Available at: www.technologyreview.com/2020/04/27/1000658/google-medical-ai-accurate-lab-real-life-clinic-covid-diabetes-retina-disease/
- 9 Haenssle, H. A., et al. (2018) Man against machine: diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists. Available at: www.sciencedirect.com/science/article/pii/S0923753419341055

Project Teams

LLOYD'S REGISTER FOUNDATION INSTITUTE FOR THE PUBLIC UNDERSTANDING OF RISK

Prof Chan Ghee Koh, Lloyd's Register Foundation Professor, Director

Prof Leonard Lee, Deputy Director

Nathaniel Tan, Senior Manager, Partnerships & Engagement

Jared Ng, Assistant Manager, Communications

Celia Leo, Communications Associate

KOREA POLICY CENTER FOR THE FOURTH INDUSTRIAL REVOLUTION

Prof So Young Kim, Director

Dr Hyeon Dae (Heidi) Rha, Senior Researcher

Dr Cornelius Kalenzi, Postdoctoral Researcher

Dr Moonjung Yim, Postdoctoral Researcher

SENSE ABOUT SCIENCE

Tracey Brown OBE, Director

Dr Hamid Khan, Programme Manager – Research Culture and Quality

Iaina Khairulzaman, Head of International Public Engagement, Training and Marketing

Joshua Gascoyne, Policy Officer

Contributors

Dr Ashraf Abdul, Affinidi

Prof Daniel Catalan, University Carlos III of Madrid

Prof Edward (Yoonjae) Choi, KAIST

Dr Sarah Cumbers, Lloyd's Register Foundation

Dr Pin Sym Foong, National University Health System

Prof Yong Jeong, KAIST

Dr Ilyoung Jung, Science & Technology Policy Institute

Dr Kyu-Hwan Jung, VUNO

Prof Steve Keevil, Guy's and St Thomas' NHS Foundation Trust

Dr Kyunghoon Kim, Korea Information Society Development Institute

Prof Tackeun Kim, Seoul National University Bundang Hospital

Prof Tze Yun Leong, NUS School of Computing

Prof Brian Lim, NUS School of Computing

Tern Poh Lim, AI Singapore

Prof Tamra Lysaght, NUS Yong Loo Lin School of Medicine

Prof Kee Yuan Ngiam, National University Health System

Prof Joon Beom Seo, University of Ulsan College of Medicine/Asan Medical Center

