



Learning with insufficient data: a multi-armed bandit perspective on covid-19 interventions

Jean Czerlinski Whitmore Ortega¹ 

Received: 8 March 2021 / Accepted: 11 October 2022 / Published online: 24 November 2022
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2022

Abstract

In February 2020, as covid-19 infections spread to more than fifty countries, public health officials needed to recommend how the public could protect themselves, balancing safety and urgency. But there was very little data since this novel virus had only been identified three months prior. How could public health officials decide with insufficient data? The multi-armed bandit problem of computer science offers adaptive decision-making procedures that can achieve both safety and urgency. These adaptive methods balance learning information (exploring) with using information (exploiting), adjusting the balance toward learning when uncertainty is high (March 1991; Kaelbling et al. 1996). Related methods are already used in adaptive clinical trials for pharmaceuticals (Pallmann et al. 2018). But we still need to develop these methods for non-pharmaceutical interventions, as I will illustrate with a case study of public mask-wearing to reduce the spread of covid-19. Public health pronouncements impact future learning.

Keywords Covid-19 · Decision-making · Public health policy · Reinforcement learning · Uncertainty · Risk

1 Introduction

On February 28, 2020, as covid-19 infections spread to more than fifty countries, the World Health Organization (WHO) increased the global risk level of covid-19 to “very high” (World Health Organization 2020). Public health officials needed to recommend how the public could protect themselves, balancing safety and urgency. But there was very little data since this novel virus had only been identified three months prior. How could public health officials decide with insufficient data?

The multi-armed bandit problem of computer science offers adaptive decision-making procedures that can achieve both safety and urgency. These adaptive

✉ Jean Czerlinski Whitmore Ortega
jeanimal@gmail.com

¹ Chicago, IL 60607, USA

methods balance *learning* information (exploring) with *using* information (exploiting), adjusting the balance toward learning when uncertainty is high (March 1991; Kaelbling et al. 1996). Related methods are already used in adaptive clinical trials for pharmaceuticals (Pallmann et al. 2018).

But we still need to develop these methods for non-pharmaceutical interventions, as I will illustrate with a case study of public mask-wearing to reduce the spread of covid-19. Public health officials recommended against face masks in February of 2020, which should have halted learning aside from randomized controlled trials. But rogue mask mandates in cities like Jena, Germany, enabled learning about the safety and efficacy of masks, leading officials to reverse course in support of face masks in June 2020. There was a randomized controlled trial of face masks, but it did not complete until September of 2021, fifteen months later. We need better methods of addressing public health emergencies than relying on rogue experiments like the one in Jena, Germany.

Let me begin by reviewing public health's historical tension between safety and urgency and then explain how adaptive methods help resolve this tension.

2 The evolution of safety vs. urgency

Public health officials have wrestled with balancing safety and urgency for decades. The thalidomide disaster of the 1950s spurred the requirement for randomized controlled trials, the gold standard of safety. But the AIDS crisis of the 1980s led to increasing use of exceptions to randomized clinical trials for desperately ill patients needing more urgent solutions.

2.1 Safety

Pharmaceutical safety grew in importance in public health. The essential rules of *randomized controlled trials* (RCTs) were established in the 1950's (Junod 2016). The rules included having both treatment and control groups, randomly assigning patients to groups, and assigning double blind (meaning neither patient nor physician knew whether the patient was in the treatment or control group). The ethics of randomized controlled trials are complicated but include a requirement of "clinical equipoise," defined as a state of genuine uncertainty in the medical community regarding the therapeutic merits of a proposed treatment (Freedman 1987).

Despite this gold standard quality, randomized controlled trials were rarely used because they were expensive and time-consuming (Junod 2016). The thalidomide disaster—which has been called the largest human-made medical disaster in history (Vargesson 2015)—changed that attitude. In the 1950's, pregnant women with morning sickness were prescribed the drug thalidomide for its antiemetic effects. Unfortunately, within a few years, doctors noticed that thalidomide was associated with a high rate of birth defects. After the dangerous side-effects of thalidomide were recognized, resistance to randomized controlled trials "suddenly melted away,"

(Junod 2016). The U.S. passed regulations requiring randomized controlled trials for pharmaceuticals in 1962.

2.2 Urgency

Then the AIDS epidemic of the 1980s presented a new problem: urgency. Young and healthy people suddenly fell ill and died. These patients demanded access to potentially life-saving new treatments *now*. They did not have years to wait for the results of randomized controlled trials (Jansen and Stryker 1993).

There was in fact a precedent for quicker access: If a patient had a life-threatening disease and had exhausted all licensed therapeutic options, the patient's doctor could request a "compassionate use" exception to treat the patient with a therapeutic agent that was already under investigation. But obtaining authorization for compassionate use was a slow, ad hoc, per-patient process. AIDS activists successfully pushed the U.S. Food and Drug Administration (FDA) to streamline the process with the 1987 "investigational new drug" (IND) regulations (Jansen and Stryker 1993). These regulations sped up AIDS patients' access to drugs. However, many AIDS activists wanted more and faster approval of drugs, so the tension between safety and urgency was not settled.

In the 2020s, the covid-19 pandemic also created urgency, and compassionate use regulations were often invoked to speed access to new treatments. For example, as early as the spring of 2020, the anti-viral remdesivir was administered to patients with severe covid-19 under the U.S. IND compassionate use regulations (Grein et al. 2020). Different countries leveraged a variety of compassionate use regulations during the covid-19 pandemic, and internationally standardized "managed access programs" have been proposed (Aliu et al. 2021).

2.3 Controversy

However, a 2020 editorial in the Journal of the American Medical Association argued that such widespread "compassionate use" risks patient safety and has many societal downsides (Kalil 2020). In particular, compassionate use harms learning: "The administration of off-label drug use, compassionate drug use, and uncontrolled studies during a pandemic also could discourage patients and clinicians from participating in RCTs, hampering any knowledge that could be gained" (Kalil 2020). That is, a desperate patient may prefer to receive a new drug under compassionate use rather than enroll in a randomized controlled trial, since a trial provides only a 50/50 chance of actually receiving the new drug rather than a placebo. (Note that urgency does not justify unscientific treatments for covid-19, like drinking bleach.)

In summary, public health historically had two extremes of decision-making procedures: the safety of randomized controlled trials or the urgency of exceptions made for compassionate use.

3 Balancing learning and using information

Computer science can shed some light on these decision-making procedures. Randomized controlled trials and compassionate use are algorithms in the “learn-by-doing” or “reinforcement learning” paradigm because we do not have the full data set up-front—we must also decide how to gather more data. In reinforcement learning, the core tradeoff is between *learning* information and *using* information, or “*exploration*” and “*exploitation*” (March 1991; Kaelbling et al. 1996). Randomized controlled trials can be seen as first focusing 100% on learning (exploring), and then, when the study is complete, focusing 100% on using (exploiting) the statistically validated treatment. In contrast, in reinforcement learning, an optimal algorithm would adaptively adjust the balance between learning information and using information based on uncertainty and trials remaining. And when there are no trials remaining, compassionate use would be optimal—doing no learning, just using (exploiting) the treatment that currently seems most promising.

Of course, public health is concerned with more than just optimality. But the tradeoff between learning and using information has already motivated the use of adaptive clinical trials of pharmaceuticals, so it is useful to understand this perspective to help us apply it to the non-pharmaceutical realm.

3.1 The two-armed bandit problem

Let me illustrate the tradeoff between learning and using information with a generic reinforcement learning problem. A “two-armed bandit” has two options that must be tried to assess success. An algorithm sequentially selects which option to try, based on past successes and failures. The goal is to design an algorithm that maximizes the total reward over time. This paper will sketch only a simple bandit algorithm. For a recent survey of the huge variety of bandit algorithms, including using multiple arms, adjusting for context, and applying behavioral or ethical constraints, see (Bouneffouf and Rish 2019).

Suppose the options are two advertisements for a video camera, one with a cat and one with a dog, and we want to know which ad garners more clicks. The advertiser assesses success by showing users an ad and measuring how many users click on the ad.

Suppose the dog ad has been shown to 500 users and garnered 5 clicks, which is a click-through rate of $5/500 = 1\%$. Meanwhile, the cat ad is new and has been shown to only 10 users and garnered zero clicks, which is a 0% click-through rate. If you had to recommend an ad at this point, which would it be? The dog ad has been more successful, so that seems to be the obvious choice. In computer science parlance, choosing the current front-runner is called the “greedy” algorithm. The dog ad is indeed optimal when there is *only one more user* to show an ad to. But what if there were one million more users? We are still uncertain about the click-through rate of the cat ad. Then showing the dog ad for the next million users *foregoes the opportunity to learn* about the true success of the cat ad.

The core insight of the bandit approach is that the “greedy” algorithm of selecting the currently best option, the “front-runner”, loses the opportunity to learn about other options. The algorithm that maximizes long-run success instead allocates some trials to learning—if there is *uncertainty* and *enough future trials* to enjoy the fruits of learning.

For example, the advertiser might show the dog ad to 80% of users and the cat ad to the other 20%. As data is gathered on click-through rates, the proportion seeing each ad shifts to favor the ad garnering more success. If the click-through rate on the dog ad continues to be higher, its proportion would climb to 90% and eventually 100%. But if the cat advertisement started getting a higher click-through rate, then the dog ad’s proportion would drop to 70, 60, 50, and eventually to 0%, which means showing only cat ads.

Showing one ad 100% of the time of course means learning has ceased. The bandit algorithms find the optimal moment to stop learning and use the fruits of learning. But before that moment, the algorithms maintain a diversity of viable options, converging on a best option only as uncertainty decreases.

3.2 A bandit perspective on RCTs and compassionate use

Now let’s apply the bandit perspective to the algorithms of randomized controlled trials and compassionate use. Randomized controlled trials use both the options equally during the entire trial—50% placebo, 50% new treatment—until the required sample size is achieved. From the bandit perspective, this would imply equal uncertainty about both treatments throughout the trial. At the end of the trial, investigators assess whether the treatment had a significant difference from the placebo: if yes, then use the treatment 100% of the time, and if not, revert to the standard of care 100% of the time. From the bandit perspective, all the exploitation of the learning takes place only after the end of the trial.

Compassionate use would be akin to an allocation of 100% new treatment, a form of greedy algorithm. Recall that the greedy algorithm is optimal if there is only one trial left, and from the perspective of a dying patient, there *is* only one trial left—their own life! But at the societal level, there are many trials left, so society would prefer to learn. This example illustrates a tension between what is optimal for a dying patient (who has one trial) and what is optimal for society (which has many trials). This tension dissipates if the patient is *not* deathly ill because the patient has a long life ahead, many future trials. Then the patient has time to benefit from the fruits of learning, too.

From the point of view of bandit algorithms, an optimal algorithm uses an *adaptive allocation* that varies based on uncertainty and trials remaining. And there actually is a form of clinical trial that leverages this idea.

3.3 Adaptive clinical trials for pharmaceuticals

Health care already has an established technique for adaptive allocation for pharmaceuticals: adaptive clinical trials. Adaptive clinical trials were allowed—but not

required—by the U.S. FDA in 2010 (Pallmann et al. 2018). Adaptive clinical trials are similar to bandit algorithms but with full statistical rigor to ensure patient safety. Adaptive trials allow investigators to review the data mid-trial to make certain types of changes to the treatment protocol. A typical change is assigning more patients to the currently more successful “arm” of the trial, similar to what bandit algorithms do. However, the adaptations are not ad-hoc: they must be carefully designed up-front by statisticians so that statistically valid conclusions can still be drawn.

With such statistical rigor, adaptive trials greatly reduce the safety versus urgency trade-off. They are just as safe as standard trials yet converge to the best treatment faster. As a recent tutorial on the topic explained: “Trials with an adaptive design are often more efficient, informative and ethical than trials with a traditional fixed design since they often make better use of resources such as time and money, and might require fewer participants” (Pallmann et al. 2018). That said, the relative merits and ethics of randomized controlled trials and adaptive clinical trials are still debated (Fillion 2019).

In any case, the urgency of the covid-19 pandemic spurred not only more compassionate use but also more adaptive clinical trials that came to conclusions more quickly than traditional randomized controlled trials. For example, an adaptive clinical trial to treat severely ill covid-19 patients with the antiviral drug remdesivir began on February 21, 2020. By April 29, 2020, just three months later, preliminary results were shared: patients who received remdesivir had a 31% faster time to recovery than those who received placebo (Beigel et al. 2020). Based on such results, the US Food and Drug Administration issued an Emergency Use Authorization for remdesivir in the treatment of covid-19 in May 2020 (Beigel et al. 2020). Adaptive clinical trials helped the world find treatments for covid-19 safely and more quickly than with traditional randomized controlled trials.

4 Challenges of non-pharmaceutical interventions

Pharmaceuticals are highly controlled. Unfortunately, making recommendations for non-pharmaceutical interventions is more challenging, both for covid-19 and other public health issues. For example, when the US Surgeon General declared in 1964 that smoking caused lung cancer, the conclusion was based on more than 7,000 studies, but none was a randomized controlled trial. Why not? Randomized controlled trials for non-pharmaceuticals tend to be logistically and ethically challenging. And when trials can be undertaken, they tend to require years. For example, the randomized controlled trial of mask-wearing for covid-19 took more than a year to organize, and results were not available until September of 2021 (Abaluck et al. 2021).

4.1 Randomized controlled trial for public mask-wearing

Why did the clinical trial for face masks take so long? For mask-wearing, the outcome measure is the infection rate of a community, which inherently requires months to measure. Investigators had to find a location with similar communities

that could participate, ultimately identifying 600 villages in rural Bangladesh. The investigators also could not mandate mask usage, so instead villages in the intervention arm were provided free masks along with encouragement and incentives to wear masks in public. It takes time to distribute masks and to broadcast public health announcements and incentives appropriate for local conditions. Because compliance was far from guaranteed, the investigators also had to measure actual mask usage. Finally, because mask-wearing is tested on healthy people, investigators had to track infection rates over five months to allow time for natural covid-19 outbreaks to develop. In the end, the first results (a preprint) were only available in September 2021. The conclusion was that a 30-percent increase in mask-wearing led to a 10 percent drop in covid-19 infection rates (Abaluck et al. 2021).

The mask-wearing trial could not easily be made into an adaptive trial because adaptive trials do not work well “if the outcome measure of interest takes so long to record that there is basically no time for the adaptive changes to come into effect before the trial end” (Pallmann et al. 2018).

The randomized controlled trial assured us of the safety of public mask-wearing—masks did not increase infection rates or have dangerous side-effects. Furthermore, masks were effective in reducing covid-19 transmission. But the first conclusions were published in September of 2021.

4.2 Uncertainty before the mask-wearing RCT

In February of 2020, public health officials were under pressure to provide recommendations to the public on protecting themselves from covid-19. At that time, using expert best judgment, several major public health agencies—the WHO, U.S. Centers for Disease Control, and European Center for Disease Prevention and Control—recommended *against* the public wearing face masks (Jingnan 2020). From the bandit point of view, using a single option halts valuable learning if uncertainty was high. So how much uncertainty was there?

Epidemiologists debated the merits of public mask-wearing in medical journals in February 2020. Some epidemiologists argued in favor of public mask-wearing as the safer option (Greenhalgh et al. 2020). Meanwhile, other epidemiologists (Martin et al. 2020) argued that masks could be harmful in several ways. First, if the public re-used dirty masks, pathogens could accumulate and increase infection risk. Second, there were possible bad behavioral side-effects. People who wore masks might not bother to use other, more effective means of reducing transmission, such as avoiding social gatherings. And even worse, these people might feel so protected by masks that they took additional risks, like visiting covid-positive people (“moral hazard”). And finally, the public might buy so many masks that healthcare professionals could not obtain an adequate supply (Martin et al. 2020).

This debate demonstrated great uncertainty and therefore a need for learning about public mask-wearing. And there was time to learn and reap the benefits of learning. Even in February 2020, we could anticipate at least a few more months of the pandemic, during which time we could be gathering data about face masks. In that case, putting 100% focus on one recommendation—in this case recommending

against masks—came at the cost of learning about the true success of wearing masks. The recommendation halted learning outside of randomized controlled trials.

4.3 Learning anyway: rogue locales and synthetic control

Or rather, the recommendation *should* have halted learning. It should have made the public wait until the randomized controlled trial was completed in September of 2021. However, some locales, such as Jena, Germany, went against the public health recommendations and required face masks in the spring of 2020 anyway. As a side-effect, these locales enabled learning about the safety and effectiveness of wearing face masks. Analyses of data from European regions (Mitze et al. 2020) and, later, US counties (Lyu and Wehby 2020) provided strong evidence that, to the surprise of health officials, wearing masks reduced covid-19 transmission. The European study summarized the effect size in lay terms:

Consider a region in which the number of COVID-19 cases increased by 10% from one day to another. This increase would have been only 6% if there had been an obligation to wear face masks. With a 10% daily increase in COVID-19, cases double within 7 days; in contrast, a 6% daily increase means cases double only within 12 days (Mitze et al. 2020).

These researchers did not have randomized controlled trials. Instead, they used a method known as “synthetic control,” which relies on existing variation—such as before and after a locale mandates masks—to make a statistically-weighted combination of groups to act as a control (Bouttell et al. 2018).

While not a gold standard, the findings of the synthetic control studies tilted the weight of evidence toward face masks. By June of 2020, public health authorities reversed course and recommended wearing face masks (Yan 2020). The randomized controlled trial for face masks continued, publishing a preprint in September 2021 (Abaluck et al. 2021), 15 months later, and confirming that the best recommendation was indeed to wear masks. Thus, the *learning enabled by cities like Jena brought public mask wearing 15 months earlier than the randomized controlled trial, saving countless lives.*

5 Public health communication

Ideally, public health policy would not rely on rogue experiments. So how could public health recommendations better balance using known information and learning additional information, doing so safely and ethically? Designing a new public health policy is beyond the scope of this essay, but I do want to leverage the bandit perspective for some potential leads in health communications, specifically not discouraging scientifically viable options and communicating uncertainty.

5.1 Not discouraging scientifically viable options

The bandit perspective emphasizes the learning opportunity cost of discouraging a scientifically viable option when uncertainty is high. While it was perfectly appropriate to declare that the balance of evidence leaned against public mask wearing in February 2020, the US Surgeon general actively discouraged and even disparaged mask-wearing. His tweet on February 29, 2020, stated: “Seriously people-STOP BUYING MASKS! They are NOT effective in preventing general public from catching #Coronavirus...” (Jingnan 2020). The tone scolds people who bought masks, even though there were well-respected epidemiologists who advocated wearing masks.

In contrast, drinking bleach to prevent or treat covid-19 was never scientifically viable and should always be discouraged. Scientists had no uncertainty about drinking bleach to cure covid-19. There was no learning cost to discouraging bleach because we had no need to learn more about bleach. Ideally, public health officials would focus discouragement where it matters most, such as not drinking bleach.

5.2 Communicating uncertainty

A more lofty communication goal would be to help the public understand all the scientifically viable options, the evidence for and against each option, and the uncertainty about each option. These are the essential inputs to a bandit algorithm and could enable people to choose options that, in aggregate, balance using information with learning new information. From a related perspective, economist Charles F. Manski has argued for “forthright communication of uncertainty in reporting of official statistics and in research that aims to inform policy” (Manski 2020). While the lay public may often seem too unsophisticated to understand uncertainties in science, there is active research on methods to successfully communicate scientific uncertainties, both verbally and visually (Hullman 2022).

Manski goes even further, arguing for actively incentivizing a diversity of policies among the scientifically plausible options when uncertainty is high. For the case of covid-19 interventions in the United States, Manski suggested, “The federal government can provide incentives to the states to encourage them to enact desirable portfolios of policies... modifying the incentives as knowledge accumulates” (Manski 2020). The intersection of incentives and diversity with safety and urgency are a topic for future research.

5.3 Adaptive methods

In summary, adaptive methods have successfully attained both safety and urgency in decision-making for pharmaceuticals under insufficient information. Non-pharmaceutical interventions are more complicated. But the covid-19 pandemic has shown the importance of urgent recommendations for these interventions. We need adaptive decision-making methods that are appropriate for issues such as mask-wearing during a pandemic while keeping a close eye on safety.

6 Changes in uncertainty

As the covid-19 pandemic has continued, our uncertainty has mostly decreased, and we have converged on a set of treatments and interventions. But as new covid-19 variants evolve, uncertainty increases, and we have to shift the balance more toward learning again. For example, in January of 2022, the CDC declared that the public needed better face masks to protect themselves from the new omicron variant (2022). This masking change is unlikely to be the last change in recommendations due to a change in the virus. We need to be ready to learn again whenever uncertainty increases again.

7 Summary

Historically, public health decision-making has had to choose between the *safety* of randomized controlled trials and the *urgency* of compassionate use that allows dying patients to try experimental treatments. But adaptive algorithms, such as those used in multi-armed bandit problems, show that public health can maintain safety while converging faster to the best option. Adaptive algorithms work by balancing learning information (“exploration”) with using information (“exploitation”), adjusting the allocation as data is gathered and uncertainty decreases. Adaptive clinical trials are already in use for pharmaceuticals and helped more quickly identify safe and effective treatments during the covid-19 pandemic. But we do not yet have adaptive decision-making for non-pharmaceuticals, such as whether the public should wear face masks. I hope the bandit algorithms’ perspective will spur the design of better public health decision-making and communication methods so we do not have to rely on rogue experiments by cities like Jena, Germany. Let us keep in mind that public health pronouncements impact future learning.

Acknowledgements I would like to thank Heather Adkins, Guido Fioretti, Dan Goldstein, Marc Harper, Jake Hofman, Donald MacKenzie, Jesus M. Ortega, Andreas Ortmann, Riccardo Rebonato, Joshua Safyan, Damian Sullivan, Oscar Wijsman, and two anonymous reviewers.

Author’s Note

Any views or opinions expressed in this piece are solely those of the author.

References

- (2022) Are cloth masks enough to protect against omicron? In: <https://Clevel.Clin.health.clevelandclinic.org/are-cloth-masks-enough-against-omicron/>. Accessed 31 Mar 2022
- Abaluck J, Kwong LH, Styczynski A et al (2021) Impact of community masking on COVID-19: a cluster-randomized trial in Bangladesh. *Science* 375:eabi9069. <https://doi.org/10.1126/science.abi9069>
- Aliu P, Sarp S, Fitzsimmons P (2021) Increasing use of compassionate use/managed access channels to obtain medicines for use in COVID-19. *Clin Pharmacol Ther* 110:26–28. <https://doi.org/10.1002/cpt.2140>

- Beigel JH, Tomashek KM, Dodd LE et al (2020) Remdesivir for the treatment of Covid-19—final report. *N Engl J Med* 383:1813–1826. <https://doi.org/10.1056/NEJMoa2007764>
- Bouneffouf D, Rish I (2019) A survey on practical applications of multi-armed and contextual bandits. *Cs Stat arXiv:1904.10040*
- Bouttell J, Craig P, Lewsey J et al (2018) Synthetic control methodology as a tool for evaluating population-level health interventions. *J Epidemiol Community Health* 72:673–678. <https://doi.org/10.1136/jech-2017-210106>
- Fillion N (2019) Clinical equipoise and adaptive clinical trials. *Topoi* 38:457–467. <https://doi.org/10.1007/s11245-018-9540-x>
- Freedman B (1987) Equipoise and the ethics of clinical research. *N Engl J Med* 317:141–145. <https://doi.org/10.1056/NEJM198707163170304>
- Greenhalgh T, Schmid MB, Czypionka T et al (2020) Face masks for the public during the covid-19 crisis. *BMJ* 369. <https://doi.org/10.1136/bmj.m1435>
- Grein J, Ohmagari N, Shin D et al (2020) Compassionate use of remdesivir for patients with severe covid-19. *N Engl J Med* 382:2327–2336. <https://doi.org/10.1056/NEJMoa2007016>
- Hullman J (2022) Advice for the government on communicating uncertainty | Statistical Modeling, Causal Inference, and Social Science. <https://statmodeling.stat.columbia.edu/2022/03/25/advice-for-the-government-on-communicating-uncertainty/>. Accessed 30 Mar 2022
- Jansen A, Stryker J (1993) Clinical research and drug regulation. In: National Research Council (US) Panel on Monitoring the Social Impact of the AIDS Epidemic (ed) *The Social Impact of AIDS in The United States*. National Academies Press (US), Washington (DC), pp 80–116
- Jingnan H (2020) Why There Are So Many Different Guidelines For Face Masks For The Public. In: NPR.org. <https://www.npr.org/sections/goatsandsoda/2020/04/10/829890635/why-there-so-many-different-guidelines-for-face-masks-for-the-public>. Accessed 3 August 2020
- Junod S (2016) FDA and clinical drug trials: a short history. In: *A quick guide to clinical trials*, 2nd edn. BioPlan Associates, Inc, Rockville, pp 29–62
- Kaelbling LP, Littman ML, Moore AW (1996) Reinforcement learning: a survey. *J Artif Intell Res* 4:237–285. <https://doi.org/10.1613/jair.301>
- Kalil AC (2020) Treating COVID-19—off-label drug use, compassionate use, and randomized clinical trials during pandemics. *JAMA* 323:1897–1898. <https://doi.org/10.1001/jama.2020.4742>
- Lyu W, Wehby GL (2020) Community use of face masks and COVID-19: evidence from a natural experiment of state mandates In The US. *Health Aff (Millwood)* 39:1419–1425. <https://doi.org/10.1377/hlthaff.2020.00818>
- Manski CF (2020) Forming COVID-19 policy under uncertainty. *J Benefit-Cost Anal* 11:341–356. <https://doi.org/10.1017/bca.2020.20>
- March JG (1991) Exploration and exploitation in organizational learning. *Organ Sci* 2:71–87. <https://doi.org/10.1287/orsc.2.1.71>
- Martin GP, Hanna E, Dingwall R (2020) Response to Greenhalgh et al.: Face masks, the precautionary principle, and evidence-informed policy. *Br Medial J* 369. <https://doi.org/10.1136/bmj.m1435>
- Mitze T, Kosfeld R, Rode J, Wälde K (2020) Unmasked! The effect of face masks on the spread of COVID-19. In: <https://voxeu.org/article/unmasked-effect-face-masks-spread-covid-19>. Accessed 14 Jul 2020
- Pallmann P, Bedding AW, Choodari-Oskooei B et al (2018) Adaptive designs in clinical trials: why use them, and how to run and report them. *BMC Med* 16:29. <https://doi.org/10.1186/s12916-018-1017-7>
- Vargesson N (2015) Thalidomide-induced teratogenesis: History and mechanisms. *Birth Defects Res Part C Embryo Today Rev* 105:140–156. <https://doi.org/10.1002/bdrc.21096>
- World Health Organization (2020) Coronavirus disease 2019 (COVID-19): situation report, 39
- Yan H (2020) Face masks in the US: Why guidance has changed so much—CNN. <https://www.cnn.com/2020/07/19/health/face-masks-us-guidance/index.html>. Accessed 3 January 2021

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

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.