

## Δεδομένα πραγματικού κόσμου: Δυνατότητες, περιορισμοί και μελλοντικές προοπτικές

Δημοσθένης Παναγιωτάκος<sup>#\*</sup> και Βίκτωρας Γκοτζαμάνης<sup>#</sup>

### Περίληψη

Η όλο και ευρύτερη χρήση του διαδικτύου, η συνεχώς αυξανόμενη χρήση έξυπνων κινητών τηλεφώνων και άλλων φορητών συσκευών με δυνατότητες καταγραφής, αποστολής και ανάλυσης δεδομένων και μια πληθώρα άλλων εφαρμογών που γίνονται διαθέσιμες με τις νέες τεχνολογίες έχουν δημιουργήσει μία πολύ σημαντική δεδομένων, τα δεδομένα πραγματικού κόσμου. Τα δεδομένα αυτά έχουν τη δυνατότητα να αξιοποιηθούν για την επιβεβαίωση ευρημάτων που προέρχονται από κλινικές μελέτες ή βασική έρευνα, αλλά και για την προσέγγιση ερωτημάτων που δε θα ήταν εφικτό να απαντηθούν με τις παραδοσιακές μεθόδους με τελικό σκοπό τη δημιουργία ιατρικών ενδείξεων και κατευθυντήριων οδηγιών. Ωστόσο, τα δεδομένα αυτά έχουν κάποιους εγγενείς περιορισμούς όπως είναι η ετερογένεια, η ασυνέπεια, η μη αντιπροσωπευτικότητα του γενικού πληθυσμού και γενικότερα η χαμηλότερη ποιότητα τους συγκριτικά με τα δεδομένα που προέρχονται από κλινικές μελέτες, λόγω του μεγάλου τους όγκου, τη μη πιστοποιημένη αξιοπιστία των μέσων καταγραφής και την ποικιλία της διαθεσιμότητας τους σε διαφορετικές ομάδες πληθυσμού. Για το λόγο αυτό απαιτείται η χρήση και περαιτέρω ανάπτυξη σύνθετων και καινοτόμων μεθόδων όπως είναι η μηχανική μάθηση και η τεχνητή νοημοσύνη για την ανάλυση τους και την εξαγωγή ασφαλών συμπερασμάτων μέσα από αυτά. Παρόλο που απαιτείται ακόμα αρκετή έρευνα για την τελειοποίηση των τεχνικών ανάλυσης, υπάρχουν ήδη παραδείγματα αξιοποίησης των δεδομένων πραγματικού κόσμου για αντιμετώπιση προκλήσεων όπως ήταν η πανδημία COVID 19 και η προσδοκία είναι πως στο μέλλον τα δεδομένα αυτά θα διαδραματίζουν όλο και σημαντικότερο ρόλο στη λήψη αποφάσεων για τη δημόσια υγεία.

---

<sup>#</sup>Τμήμα Επιστήμης Διαιτολογίας και Διατροφής, Σχολή Επιστημών Υγείας & Αγωγής, Χαροκόπειο Πανεπιστήμιο Αθηνών  
<sup>\*</sup>Σχολή Επιστημών Υγείας, Πανεπιστήμιο της Καμπέρα, Αυστραλία

## **Εισαγωγή**

Τα δεδομένα πραγματικού κόσμου στην ιατρική είναι τα δεδομένα εκείνα που σχετίζονται με την υγεία των ασθενών και/ή την παροχή υπηρεσιών υγείας και συλλέγονται από διάφορες πηγές (1). Η ευρεία χρήση του διαδικτύου, των μέσων κοινωνικής δικτύωσης, των φορητών ηλεκτρονικών συσκευών, η δημιουργία μητρώων ασθενειών, τα ηλεκτρονικά αρχεία καταγραφής ιατρικών πληροφοριών, ηλεκτρονικές υπηρεσίες υγείας και άλλες υπηρεσίες που βασίζονται στις νέες τεχνολογίες σε συνδυασμό με την αυξημένη δυνατότητα αποθήκευσης ψηφιακής πληροφορίας έχουν καταστήσει δυνατή τη δημιουργία αλλά και ευρεία διαθεσιμότητα δεδομένων πραγματικού κόσμου (2).

Η αυξανόμενη προσβασιμότητα σε δεδομένα πραγματικού κόσμου και η ταχεία ανάπτυξη τεχνικών μηχανικής μάθησης και της τεχνητής νοημοσύνης, σε συνδυασμό με τα αυξανόμενα κόστη και τους εγνωσμένους περιορισμούς των παραδοσιακών κλινικών δοκιμών έχουν οξύνει το ενδιαφέρον για χρήση δεδομένων πραγματικού κόσμου με σκοπό την ενίσχυση της αποτελεσματικότητας της κλινικής έρευνας και τη γεφύρωση του χάσματος μεταξύ των ευρημάτων από τις κλινικές έρευνες και της καθ'ημέρα πράξης. Επί παραδείγματι, κατά τη διάρκεια της πανδημίας COVID-19, δεδομένα πραγματικού κόσμου χρησιμοποιήθηκαν για τη δημιουργία ιατρικών ενδείξεων σχετικά με την αποτελεσματικότητα του εμβολιασμού έναντι του Sars-Cov-2 (3,4,5), για τη μοντελοποίηση στρατηγικών περιορισμού της νόσου (6), για την καταγραφή των περιστατικών COVID-19 και γρίπης (7), για τη μελέτη του ψυχολογικού αντίκτυπου των περιορισμών (8) στις μετακινήσεις και για τη διαμόρφωση πολιτικών αντιμετώπισης.

## **Χαρακτηριστικά, είδη και εφαρμογές των δεδομένων πραγματικού κόσμου**

Τα δεδομένα πραγματικού κόσμου έχουν κάποια συγκεκριμένα χαρακτηριστικά συγκριτικά με τα δεδομένα που προέρχονται από τυχαίοποιημένες μελέτες σε ελεγχόμενα περιβάλλοντα. Αρχικά, τα δεδομένα πραγματικού κόσμου είναι προϊόντα παρατήρησης

χωρίς εξωτερική παρέμβαση. Δεύτερον, δεν είναι δομημένα και συχνά είναι ασυνεπή λόγω της ετερογένειας στην καταγραφή τους ανάλογα με την πηγή τους. Τρίτον, δημιουργούνται με ταχείς ρυθμούς με αποτέλεσμα να είναι ογκώδη και δυναμικά. Τέταρτον, μπορεί να είναι ελλιπή ως προς βασικές εκβάσεις καθώς η αρχική κατάγραφή τους δε γίνεται με βάση ένα συγκεκριμένο σκοπό. Για παράδειγμα, δεδομένα από μητρώα καταγραφής ασθενειών έχουν περιορισμένη πληροφορία σχετικά με την μετέπειτα παρακολούθηση. Πέμπτον, τα δεδομένα πραγματικού κόσμου είναι πιο πιθανό να υπόκεινται σε σφάλματα (τυχαία ή συστηματικά). Παραδείγματος χάριν, δεδομένα από φορητές ηλεκτρονικές συσκευές όπως κινητά ή έξυπνα ρολόγια μπορεί να υπόκεινται σε σφάλμα επιλογής αφού το δείγμα από το οποίο προέρχονται δεν είναι αντιπροσωπευτικό του γενικού πληθυσμού. Εν ολίγοις, τα δεδομένα πραγματικού κόσμου είναι ελλιπή, ετερογενή και επιρρεπή σε σφάλματα διαφόρων τύπων. Μια συστηματική ανασκόπηση έδειξε ότι η ποιότητα τέτοιων δεδομένων είναι τόσο ετερογενής που καθιστά ιδιαίτερα δυσχερή ακόμα και την καταγραφή της σε μεγάλη κλίμακα. Η χαμηλή ποιότητα λοιπόν, των δεδομένων πραγματικού κόσμου είναι αναγνωρισμένη (9,10,11,12) αλλά ο τρόπος για τη βελτιστοποίησή της ακόμα αναζητείται (13,14,15).

Όπως αναφέρθηκε, υπάρχουν αρκετές πηγές δεδομένων πραγματικού κόσμου. Θα αναφερθούμε ενδεικτικά σε ορισμένες από αυτές για να αναδείξουμε την ποικιλία και πώς μπορούν να αξιοποιηθούν: Δεδομένα από ηλεκτρονικά αρχεία καταγραφής ιατρικών πληροφοριών, δεδομένα από μητρώα ασθενειών, δεδομένα από ασφαλιστικά ταμεία, δεδομένα από καταγραφή των ίδιων των ασθενών και δεδομένα από φορητές ηλεκτρονικές συσκευές.

Ηλεκτρονικά αρχεία καταγραφής ιατρικών πληροφοριών χρησιμοποιούνται στην καθημέρα πράξη σε κλινικές, νοσοκομεία και άλλες δομές παροχής υπηρεσιών υγείας. Τα αρχεία αυτά αποτελούν ένα χαρακτηριστικό παράδειγμα δεδομένων πραγματικού κόσμου: Είναι ετερογενή, δυναμικά και απαιτείται ενταντική προσπάθεια για την προπαρασκευή τους προτού να μπορούν να χρησιμοποιηθούν για ανάλυση (16). Τα αρχεία αυτά έχουν δημιουργήσει νέες ευκαιρίες για την ανάπτυξη κλινικών προσεγγίσεων,

την ανάδειξη μοτίβων, τη βελτίωση του περιεγχειρητικού σχεδιασμού, την ακριβέστερη και ταχύτερη διάγνωση, την καλύτερη πρόγνωση καθώς και την επικύρωση και την αναπαραγωγή ευρημάτων από κλινικές δοκιμές (17-30), ιδιαίτερος όταν συνδυάζονται και με άλλες πηγές πληροφοριών και αναλύονται με τεχνικές μηχανικής μάθησης.

Υπάρχουν διάφοροι τύποι μητρώων καταγραφής. Μητρώα προϊόντων καταγράφουν ασθενείς που έχουν λάβει κάποιο ιατροφαρμακευτικό σκεύασμα ή συσκευή, μητρώα υπηρεσιών υγείας περιλαμβάνουν ασθενείς που έχουν υποβληθεί σε μια συγκεκριμένη χειρουργική επέμβαση ή έχουν νοσηλευτεί, μητρώα ασθενειών περιλαμβάνουν ασθενείς που έχουν διαγνωστεί με μία συγκεκριμένη πάθηση. Δεδομένα από τέτοιου είδους μητρώα συμβάλλουν στην αναγνώριση και κοινοποίηση των βέλτιστων κλινικών πρακτικών, βελτιώνουν την ακρίβεια υπολογισμών και παρέχουν πολύτιμη πληροφορία για την λήψη ρυθμιστικών αποφάσεων (31-34). Ειδική μνεία πρέπει να γίνει στις σπάνιες παθήσεις για τις οποίες τα δεδομένα από κλινικές μελέτες βασίζονται σε μικρά δείγματα και παρουσιάζουν υψηλή ετερογένεια, ενώ τα μητρώα καταγραφής τους παρέχουν πολύτιμη πληροφορία για την κατανόηση της πορείας της νόσου και μπορούν να καθοδηγήσουν τη διενέργεια επιβεβαιωτικών κλινικών μελετών και περαιτέρω βασικής έρευνας για την ανάπτυξη θεραπειών (35,36).

Τα δεδομένα από ασφαλιστικά ταμεία αν και συλλέγονται με βασικό σκοπό τη χορήγηση αποζημιώσεων και την κάλυψη εξόδων, έχουν χρησιμοποιηθεί για να βοηθήσουν στην κατανόηση της συμπεριφοράς των ασθενών και των συνταγογραφούντων και τις αλληλεπιδράσεις τους, στον υπολογισμό του επιπολασμού ασθενειών, στην καταγραφή της χρήσης φαρμάκων και αλληλεπιδράσεων μεταξύ θεραπειών καθώς και στην επιβεβαίωση και επικύρωση ευρημάτων από κλινικές μελέτες (37-44). Ωστόσο, τέτοια δεδομένα ενέχουν τον κίνδυνο να είναι ανακριβή στα πλαίσια απάτης για οικονομικούς λόγους. Με τη χρήση, όμως, ελεγκτικών μηχανισμών και σωστών τεχνικών ανάλυσης το πρόβλημα αυτό μπορεί να αντιμετωπιστεί σε σημαντικό βαθμό (45-48).

Τα δεδομένα από καταγραφή των ίδιων των ασθενών προέρχονται παρέχουν πληροφορίες σχετικά με την αποτελεσματικότητα παρεμβάσεων, την παρακολούθηση

συμπτωμάτων και την σχέση μεταξύ διαφόρων εκθέσεων και εκβάσεων (49-52). Τέτοιου είδους δεδομένα μπορεί να υπόκεινται σε σφάλμα λόγω πλημμελούς ανάκλησης από τον ασθενή ή μεγάλης διακύμανσης μεταξύ των ασθενών.

Τέλος, οι διάφορες φορητές συσκευές έχουν τη δυνατότητα συνεχούς ροής παροχής δεδομένων. Όταν συνδυάζονται με δεδομένα από άλλες πηγές δημιουργούν τη δυνατότητα για διεξαγωγή ευρείων μελετών που διαφορετικά θα ήταν αδύνατες σε συνθήκες κλινικής μελέτης (53). Οι συσκευές αυτές δημιουργούν πολύ μεγάλες ποσότητες δεδομένων και απαιτείται πρόοδος στις τεχνολογίες αποθήκευσης δεδομένων, στην ενεργειακή απόδοση μπαταριών και στην δυνατότητα ανάλυσης δεδομένων σε πραγματικό χρόνο για την πλήρη αξιοποίησή τους.

### **Χρήση και ανάλυση δεδομένων πραγματικού κόσμου**

Υπάρχει μια ευρεία γκάμα μεθόδων για τη χρήση των δεδομένων πραγματικού κόσμου. Θα αναφερθούμε ενδεικτικά σε κάποιες από αυτές όπως: Πρακτικές κλινικές δοκιμές, δοκιμές προσομοίωσης και εφαρμογές μηχανικής μάθησης και τεχνητής νοημοσύνης.

Οι πρακτικές κλινικές μελέτες είναι σχεδιασμένες για να ελέγχουν την αποτελεσματικότητα μιας παρέμβασης σε κλινικές συνθήκες πραγματικού κόσμου. Οι μελέτες αυτές εκμεταλλεύονται την τα ενοποιημένα συστήματα υγείας και μπορούν να αξιοποιήσουν δεδομένα από ηλεκτρονικά αρχεία καταγραφής, από ασφαλιστικά ταμεία, από καταγραφή από τους ίδιους τους ασθενείς κλπ. Λόγω των περιορισμών των δεδομένων πραγματικού κόσμου αναπτύσσονται νέες κατευθυντήριες οδηγίες και μεθοδολογίες για να ελαχιστοποιήσουν τα συστηματικά σφάλματα κατά τη διάρκεια δημιουργίας ιατρικών ενδείξεων μέσα από δεδομένα πραγματικού κόσμου (54,55). Το κλινικό ερώτημα που τίθεται σε τέτοιου είδους μελέτες είναι αν η παρέμβαση δουλεύει σε περιβάλλον καθημερινής πρακτικής και σχεδιάζονται για να μεγιστοποιούν την εφαρμοσιμότητα και την δυνατότητα γενίκευσης της παρέμβασης. Διάφορες εκβάσεις μπορούν να εκτιμηθούν στις μελέτες αυτές, έχοντας όμως κατά κύριο λόγο τον ασθενή στο επίκεντρο, εν αντιθέσει με τα συνηθισμένα μετρήσιμα συμπτώματα ή δείκτες που χρησιμοποιούνται στις επεξηγηματικές μελέτες. Επί παραδείγματι, η μελέτη ADAPTABLE (56,57) είναι μία

πρακτική κλινική μελέτη και είναι η πρώτη μελέτη μεγάλης κλίμακας που χρησιμοποίησε ηλεκτρονικά ιατρικά αρχεία καταγραφής στις Ηνωμένες Πολιτείες. Αξιοποίησε τα δεδομένα αυτά για να αναγνωρίσει περίπου 450.000 ασθενείς με εγκατεστημένη αθηροσκληρωτική καρδιαγγειακή νόσο ως πιθανούς συμμετέχοντες και τελικά να εντάξει 15.000 από αυτούς σε 40 κλινικά κέντρα και να τους τυχαιοποιήσει σε δύο ομάδες λήψης διαφορετικών δόσεων ασπιρίνης. Η παρακολούθηση έγινε ηλεκτρονικά με καταγραφή από τους ίδιους τους ασθενείς κάθε 3 με 6 μήνες, με διάμεση διάρκεια παρακολούθησης τους 26,2 μήνες, για να αποφασιστεί η βέλτιστη δόση ασπιρίνης ασθενείς με καρδιαγγειακή νόσο με βασικό σύνθετο καταληκτικό σημείο τη θνησιμότητα ανεξαρτήτως αιτίας, την νοσηλεία για μη θανατηφόρο έμφραγμα μυοκαρδίου, ή τη νοσηλεία για μη θανατηφόρο αγγειακό εγκεφαλικό επεισόδιο. Το κόστος της μελέτης ADAPTABLE υπολογίζεται πως ήταν μόνο το 1/5 με 1/2 μιας τυχαιοποιημένης κλινικής μελέτης τέτοιας κλίμακας.

Οι μελέτες προσομοίωσης προκύπτουν από την εφαρμογή του σχεδιασμού και των αρχών της ανάλυσης από μία επιλεγμένη τυχαιοποιημένη τυχαιοποιημένη κλινική δοκιμή για την ανάλυση δεδομένων που προέρχονται από παρατήρηση (58). Καθορίζοντας σαφώς τα κριτήρια ένταξης και αποκλεισμού, θεραπείες, εκβάσεις, την περίοδο παρακολούθησης και τη στατιστική ανάλυση μιας επιλεγμένης κλινικής δοκιμής, είναι εφικτό να εξαχθούν έγκυρες κλινικές πληροφορίες για μία παρέμβαση από δεδομένα πραγματικού κόσμου. Οι μελέτες προσομοίωσης είναι ένα πολύτιμο εργαλείο, ιδιαίτερα όταν δεν υπάρχουν διαθέσιμα συγκριτικά δεδομένα από τυχαιοποιημένες κλινικές μελέτες. Επί παραδείγματι, μελέτες προσομοίωσης αξιοποιήθηκαν για να αξιολογήσουν την αποτελεσματικότητα του εμβολιασμού έναντι του Sars-Cov-2 στην πρόληψη της λοίμωξης και της θνητότητας της νόσου σε φυλετικά διαφορετικούς πληθυσμούς ηλικιωμένων, συγκρίνοντας πρόσφατα εμβολιασμένα άτομα με αντίστοιχισμένους μη εμβολιασμένους μάρτυρες χρησιμοποιώντας δεδομένα από το τμήμα υποθέσεων υγείας βετεράνων των Ηνωμένων Πολιτειών (59). Οι μελέτες αυτές διεξήχθησαν με καθορισμένα κριτήρια ένταξης και αποκλεισμού, με ένταξη αντιστοιχισμένων μαρτύρων ως προς τα κλινικά σημαντικά χαρακτηριστικά και προσεκτικό έλεγχο για όλους τους πιθανούς

συγχυτικούς παράγοντες. Μελέτες προσομοίωσης έχουν αξιοποιηθεί επίσης για να αξιολογήσουν την επίδραση του προληπτικού ελέγχου για καρκίνο παχέος εντέρου στην επίπτωση του καρκίνου αυτού σε διάστημα παρακολούθησης 8 ετών (60) καθώς και στον κίνδυνο λοίμωξης ουροποιητικού σε διαβητικούς ασθενείς (61).

Δεδομένα πραγματικού κόσμου μπορούν να χρησιμοποιηθούν για την ανεύρεση μαρτύρων και ομάδων αναφοράς για κλινικές δοκιμές, μετά από αξιολόγηση της ποιότητας και της καταλληλότητας τους και την εφαρμογή των κατάλληλων στατιστικών προσεγγίσεων για την ανάλυσή τους (62). Ο έλεγχος για συστηματικό σφάλμα επιλογής είναι καθοριστικός για την εγκυρότητα αυτής της προσέγγισης λόγω της ελλειψης τυχαιοποίησης και πιθανώς διαλαθουσών διαφορών κατά την αρχική κατάσταση μεταξύ θεραπευομένων και μαρτύρων. Τα δεδομένα πραγματικού κόσμου δίνουν επίσης τη δυνατότητα για μελέτη σπάνιων παθήσεων και συμβαμάτων. Τα ανωτέρω όμως αναδεικνύουν την ανάγκη για βελτίωση της ποιότητας των δεδομένων πραγματικού κόσμου καθώς και την τυποποίησή τους με σκοπό τη μείωση της ετερογένειας και την διευκόλυνση της ανάλυσής τους (63-65).

Για την ανάλυση δεδομένων πραγματικού κόσμου είναι απαραίτητη η χρησιμοποίηση στατιστικών μοντέλων και επαγωγικών προσεγγίσεων, από τα οποία θα προκύπτουν πιθανές αιτιολογικές συσχετίσεις και θα ελέγχονται ή επικυρώνονται ερευνητικές υποθέσεις με τελικό σκοπό την δημιουργία ιατρικών ενδείξεων στις οποίες θα βασίζονται οι διάφοροι ρυθμιστικοί μηχανισμοί για την έκδοση οδηγιών, όπως γίνεται και με τις τυχαιοποιημένες κλινικές μελέτες. Ωστόσο, η χρησιμοποίηση των πρακτικών κλινικών δοκιμών και των μελετών προσομοίωσης για ανεύρεση αιτιολογικών συσχετίσεων δημιουργεί την ανάγκη για χρησιμοποίηση πιο καινοτόμων μεθόδων ανάλυσης από τις παραδοσιακές στατιστικές μεθόδους για έλεγχο πιθανών συμπαραγόντων, με τις οποίους θα αυξάνεται η ισχύς των δεδομένων πραγματικού κόσμου, παρά τους εγγενείς περιορισμούς τους (66-69).

Οι τεχνικές μηχανικής μάθησης γίνονται όλο και περισσότερο δημοφιλείς και είναι ισχυρά εργαλεία για προβλεπτικά μοντέλα. Ένας λόγος για τη δημοτικότητα τους είναι ότι οι

μέθοδοι αυτές είναι ικανές να αξιοποιούν ογκώδη και αδόμητα δεδομένα διαφορετικών τύπων, χωρίς να απαιτούν εκ των προτέρων υποθέσεις ως προς την πιθανή κατανομή τους. Για παράδειγμα, η βαθιά μάθηση (deep learning) μπορεί να εκπαιδευτεί σε αφηρημένες αναπαραστάσεις μεγάλων, πολύπλοκων και αδόμητων δεδομένων. Η επεξεργασία φυσικής γλώσσας και άλλες μέθοδοι ενσωμάτωσης μπορούν να χρησιμοποιηθούν για την επεξεργασία κειμένων και κλινικών σημειώσεων από ηλεκτρονικά αρχεία καταγραφής ιατρικών πληροφοριών και να τα μετασχηματίσουν σε μορφές οι οποίες μπορούν να χρησιμοποιηθούν για περαιτέρω ανάλυση και εκπαίδευση μοντέλων. Επιπλέον, νέες και πιο ισχυρές τεχνικές μηχανικής μάθησης αναπτύσσονται ραγδαία λόγω της υψηλής ζήτησης και του μεγάλου αριθμού ερευνητών που απασχολούνται στο πεδίο. Επίσης, υπάρχουν πολλά προγράμματα με ανοιχτά διαθέσιμο κώδικα και βιβλιοθήκες λογισμικών που διευκολύνουν την εφαρμογή των τεχνικών αυτών. Πράγματι την τελευταία δεκαετία η μηχανική μάθηση έχει βρει αρκετές εφαρμογές στην ανάλυση δεδομένων πραγματικού κόσμου, υπερτερώντας σημαντικά συγκριτικά με τις συμβατικές μεθόδους (70-78). Επί παραδείγματι, η μηχανική μάθηση χρησιμοποιείται στην ιατρική πληροφορική για τη δημιουργία ιατρικών ενδείξεων και τη διαμόρφωση εξατομικευμένων θεραπευτικών προσεγγίσεων (79-83). Χρησιμοποιήθηκε με επιτυχία σε δεδομένα πραγματικού κόσμου για τη διάρκεια της πανδημίας COVID-19 για την κατανόηση της νόσου και την αξιολόγηση μεθόδων πρόληψης και θεραπείας (84-88). Πρέπει να σημειωθεί ότι στην παρούσα φάση οι τεχνικές μηχανικής μάθησης χρησιμοποιούνται κυρίως για πρόβλεψη, ταξινόμηση (π.χ. διάγνωση), επιλογή σημαντικών παραμέτρων (π.χ. ανάδειξη βιοδεικτών), οπτικοποίηση αποτελεσμάτων και όχι για τη δημιουργία ιατρικών ενδείξεων. Αυτό όμως είναι κάτι που μπορεί να αλλάξει σύντομα αφού τα αποτελέσματα από μοντέλα μηχανικής μάθησης αξιολογούνται συνεχώς και εφόσον η αποτελεσματικότητά τους αποδειχτεί και στην κλινική πράξη θα διευρυνθεί και η χρησιμοποίησή τους ακόμα περισσότερο (89-92).

Ο συνδυασμός μηχανικής μάθησης και εξειδικευμένης γνώσης από στατιστικά μοντέλα είναι πιο αποτελεσματικός για τη δημιουργία ιατρικών ενδείξεων και την ανάδειξη αιτιολογικών σχέσεων. Μία πρόσφατη εξέλιξη είναι πράγματι σε αυτή την κατεύθυνση, η



αξιοποίηση των προόδων σε ημι-παραμετρικές και εμπειρικές θεωρίες επεξεργασίας σε συνδυασμό με την ενσωμάτωση των πλεονεκτημάτων της μηχανικής μάθησης στην συγκριτική αποτελεσματικότητα χρησιμοποιώντας δεδομένα πραγματικού κόσμου. Ένα γνωστό πλαίσιο είναι η στοχευμένη μάθηση (93-95) που έχει χρησιμοποιηθεί με επιτυχία στην αξιολόγηση της αποτελεσματικότητας θεραπειών για την COVID-19 (96).

Ανεξάρτητα από την πηγή δεδομένων πραγματικού κόσμου στην οποία εστιάζεται το ενδιαφέρον, είναι εξαιρετικά σημαντικό να προσδιοριστεί η δυνατότητα γενίκευσης των ευρημάτων που προκύπτουν. Σε κάθε άλλη περίπτωση υπολογισμοί και προβλέψεις μπορεί να είναι παραπλανητικές και επικίνδυνες. Η πληροφορία η οποία περιέχεται σε δεδομένα πραγματικού κόσμου μπορεί να μην επαρκεί για να επικυρωθεί η δυνατότητα γενίκευσης των ευρημάτων που προκύπτουν από αυτά. Σε αυτή την περίπτωση, οι ερευνητές θα πρέπει να είναι επιφυλακτικοί στο να διατυπώνουν επίφοβες γενικεύσεις για τις οποίες δεν είναι απόλυτα σίγουροι.

### **Δυσκολίες, Ευκαιρίες και ζητήματα ηθικής**

Διάφορες δυσκολίες, από τη συλλογή των δεδομένων μέχρι τον έλεγχο της ποιότητάς τους και τη λήψη αποφάσεων, υπάρχουν ακόμα σε όλα τα στάδια του κύκλου ζωής των δεδομένων πραγματικού κόσμου, παρ'όλο τον ενθουσιασμό γύρω από τις δυνατότητες που προσφέρουν. Θα αναφέρουμε μερικές ενδεικτικά για τις οποίες υπάρχουν προοπτικές βελτίωσης και άλλες που απαιτείται μεγαλύτερη προσπάθεια για να ξεπεραστούν.

**Η ποιότητα των δεδομένων:** Τα δεδομένα πραγματικού κόσμου συχνά χρησιμοποιούνται για διαφορετικούς σκοπούς από αυτούς για τους οποίους συλλέχθηκαν αρχικά και κατά συνέπεια μπορεί να είναι ελλιπή ως προς σημαντικές πληροφορίες που απαιτούνται για τη δημιουργία πειστικών ενδείξεων που μπορούν να αποτελέσουν τη βάση για έκδοση οδηγιών και ρυθμίσεων. Επιπλέον, τα δεδομένα πραγματικού κόσμου είναι ακατάστατα, ετερογενή και επιρρεπή σε διάφορα σφάλματα καταγραφής. Όλα αυτά συμβάλλουν στη χαμηλότερη ποιότητα των δεδομένων πραγματικού κόσμου συγκριτικά με δεδομένα από τυχαιοποιημένες κλινικές δοκιμές. Κατά συνέπεια, η ακρίβεια των

αποτελεσμάτων που βασίζονται σε δεδομένα πραγματικού κόσμου επηρεάζεται αρνητικά και μπορεί να εξαχθούν ανακριβή και παραπλανητικά συμπεράσματα. Αν και τα παραπάνω δεν αποκλείουν τη χρήση δεδομένων πραγματικού κόσμου από τη διαδικασία δημιουργίας ιατρικών ενδείξεων και λήψης αποφάσεων, τα προβλήματα που σχετίζονται με την ποιότητα των δεδομένων πρέπει να καταγράφονται πολύ αναλυτικά και να γίνονται όλες οι εφικτές προσπάθειες για να αντιμετωπιστούν. Αν ένα πρόβλημα αναγνωρίζεται κατά το στάδιο της προ-επεξεργασίας, πρέπει να γίνονται όλες οι απαραίτητες προσπάθειες διόρθωσης κατά το στάδιο της ανάλυσης αλλιώς η ερμηνεία των αποτελεσμάτων πρέπει να γίνεται με πολύ μεγάλη επιφύλαξη. Η έγκαιρη ανάμιξη των οργανισμών και των δομών που θα αξιοποιήσουν τα αποτελέσματα αυτά (ρυθμιστικοί μηχανισμοί, ερευνητικά ιδρύματα κλπ) πρέπει να ενθαρρύνεται ούτως ώστε να συμμετέχουν στην εξασφάλιση των καλύτερων δυνατών προδιαγραφών ποιότητας των δεδομένων και να αποφεύγονται πιθανοί κίνδυνοι εν τη γενέσει τους.

***Αποτελεσματικές και πρακτικές τεχνικές μηχανικής μάθησης και στατιστικές αναλύσεις:*** Η γρήγορη ανάπτυξη των ψηφιακών ιατρικών δεδομένων και το γεγονός πως οι εργαζόμενοι και οι επενδύσεις συρρέουν στο πεδίο οδηγούν την ταχεία ανάπτυξη και υιοθέτηση μοντέρνων στατιστικών μεθόδων και αλγορίθμων μηχανικής μάθησης για την ανάλυση των δεδομένων. Η διαθεσιμότητα προγραμμάτων ελεύθερου κώδικα διευκολύνουν την εφαρμογή των μεθόδων αυτών στην πράξη. Από την άλλη, η ετερογένεια, η έλλειψη δομής και ο θόρυβος των δεδομένων πραγματικού κόσμου μπορεί να οδηγήσει στην αποτυχία των υπαρχουσών τεχνικών και κατά συνέπεια απαιτείται η ανάπτυξη νέων μεθόδων που θα στοχεύουν συγκεκριμένα στην αξιοποίηση τέτοιων δεδομένων. Επιπλέον, η διαθεσιμότητα ελεύθερων προγραμμάτων ανοιχτού κώδικα και η διευκόλυνση που προσφέρουν, αν και προσφέρονται με καλές προθέσεις, αυξάνουν την πιθανότητα οι τεχνικές αυτές να καταλήξουν να χρησιμοποιηθούν από ανθρώπους που δεν έχουν την κατάλληλη εκπαίδευση γύρω από αυτές και δεν κατανοούν πλήρως τις αρχές που τις διέπουν. Επιπροσθέτως, για τη διατήρηση του επιστημονικής αριότητας κατά τη δημιουργία ιατρικών ενδείξεων από δεδομένα πραγματικού κόσμου, τα αποτελέσματα από στατιστικές μεθόδους και τεχνικές μηχανικής μάθησης πρέπει να

επιβεβαιώνονται και να επικυρώνονται ιατρικά μέσω της ήδη υπάρχουσας εξειδικευμένης γνώσης ή με τη διαξαγωγή μελετών που θα αναπαράγουν τα αποτελέσματα προτού μπορέσουν να χρησιμοποιηθούν για τη λήψη αποφάσεων (97).

**Ικανότητα επεξήγησης και ερμηνείας:** Οι σύγχρονες τεχνικές μηχανικές μάθησης ομοιάζουν με ένα μαύρο κουτί και λείπει η κατανόηση των σχέσεων μεταξύ των εισαγόμενων και των εξαγόμενων δεδομένων και των αιτιολογικών σχέσεων. Η επιλογή μοντέλου, οι αρχικές ρυθμίσεις του και η βελτιστοποίηση της απόδοσής του συχνά πραγματοποιούνται μέσα από διαδικασίες δοκιμής και λάθους, χωρίς την καθοδήγηση ειδικών στον τομέα. Αυτό είναι αντίθετο με τις συνήθεις διαδικασίες στον ιατρικό τομέα όπου ο τρόπος ερμηνείας είναι καθοριστικός για το χτίσιμο σχέσεων εμπιστοσύνης μεταξύ ιατρού και ασθενούς και επιπλέον οι γιατροί είναι αρκετά απίθανο να χρησιμοποιήσουν κάποια τεχνολογία που δεν κατανοούν πλήρως. Παρόλο που ήδη έχει ξεκινήσει ελπιδοφόρα έρευνα στο πεδίο (98-103), απαιτούνται ακόμα αρκετά βήματα προόδου.

**Επαναληψιμότητα και αναπαραγωγιμότητα:** Η αναπαραγωγιμότητα και επαναληψιμότητα είναι βασικές αρχές στην επιστημονική έρευνα. Η αναπαραγωγιμότητα είναι η ικανότητα παραγωγής των ίδιων ευρημάτων από την ανάλυση των δεδομένων από ανεξάρτητο ερευνητή, Η επαναληψιμότητα είναι η ικανότητα επανάληψης των ευρημάτων σε διαφορετικά δείγματα. Αν μια αναλυτική διαδικασία δεν είναι αρκετά σαφής και τα αποτελέσματά της δεν μπορούν να επαναληφθούν και να αναπαραχθούν από ανεξάρτητους ερευνητές, η επιστημονική της αρτιότητα τίθεται υπό αμφισβήτηση και τα συμπεράσματα που προκύπτουν είναι επισφαλής (104-107). Η επικύρωση των αποτελεσμάτων, η αναπαραγωγιμότητα και η επαναληψιμότητα μπορεί να είναι δύσκολο να επιτευχθούν δεδομένης της αταξίας των δεδομένων πραγματικού κόσμου αλλά πρέπει να εξασφαλίζονται ιδιαίτερα όταν τα αποτελέσματα πρόκειται να χρησιμοποιηθούν για τη λήψη αποφάσεων που μπορεί να επηρεάσουν τη ζωή εκατομμυρίων ανθρώπων. Η μη αναπαραγωγιμότητα μπορεί να αποφευχθεί μέσα από την κοινοποίηση των ακατέργαστων και επεξεργασμένων δεδομένων και του κώδικα που χρησιμοποιήθηκε,

δεδομένου ότι δεν τίθεται σε κίνδυνο η ιδιωτικότητα των δεδομένων. Έχοντας υπόψιν ότι τα δεδομένα πραγματικού κόσμου δεν παράγονται από τυχαιοποιημένες μελέτες και κάθε βάση δεδομένων μπορεί να έχει τα δικά της ιδιαίτερα χαρακτηριστικά, η επαναληψιμότητα μπορεί κατά περίπτωση να είναι δύσκολη ως και αδύνατη. Ωστόσο, η λεπτομερής καταγραφή των χαρακτηριστικών των δεδομένων και της προπαρασκευής τους καθώς και η συμμόρφωση με τις αρχές της ανοικτής επιστήμης είναι καθοριστικές για να μπορού να επαναληφθούν τα ευρήματα από δεδομένα πραγματικού κόσμου, με την προϋπόθεση ότι προέρχονται από τον ίδιο πληθυσμό.

Οι ανωτέρω δυσκολίες δεν είναι ανεξάρτητες αλλά σε σχέση μεταξύ τους, Η ποιότητα των δεδομένων επηρεάζει την απόδοση των στατιστικών μοντέλων και των τεχνικών μηχανικής μάθησης. Οι πηγές των δεδομένων και η προ-επεξεργασία των δεδομένων συνδέεται με την αναπαραγωγιμότητα και την επαναληψιμότητα. Ο τρόπος ανάλυσης και το στατικά μοντέλα που θα χρησιμοποιηθούν έχουν αντίκτυπο στην αναπαραγωγιμότητα και την επαναληψιμότητα.

### ***Ζητήματα ηθικής***

Η χρήση δεδομένων πραγματικού κόσμου εγείρει μια σειρά από ζητήματα ηθικής. Η ποικιλία των πηγών, των χαρακτηριστικών και κυρίως των τρόπων αξιοποίησης των δεδομένων πραγματικού κόσμου αντανάκλαται και στην πολυπλοκότητα των ζητημάτων ηθικής που σχετίζονται με αυτά. Οι αξίες που πρέπει να λαμβάνονται υπόψιν σε μια μελέτη πραγματικού κόσμου διαφοροποιούνται ανάλογα με τη σκοπιά των διάφορων εμπλεκόμενων(108).

***Ιδιωτικότητα:*** Ηθικά ζητήματα προκύπτουν όταν τίθεται υπό συζήτηση η χρήση δεδομένων πραγματικού κόσμου και η ιδιωτικότητα είναι ένα συχνό θέμα που ανακύπτει. Οι πληροφορίες στα δεδομένα πραγματικού κόσμου είναι συχνά ευαίσθητες, όπως το ιατρικό ιστορικό, κάποια νόσηση, οικονομική κατάσταση, κοινωνικές συμπεριφορές και

άλλα. Η ιδιωτικότητα μπορεί να τεθεί υπό κίνδυνο όταν διαφορετικές βάσεις δεδομένων συνδέονται, μια κοινή πρακτική κατά τη χρήση δεδομένων πραγματικού κόσμου. Η ιδιωτικότητα μπορεί να προστατευτεί με την ανωνυμοποίηση των δεδομένων καθώς και με την επαρκή ενημέρωση και δήλωση συγκατάθεσης των συμμετεχόντων (109). Ωστόσο με τις μεθόδους αυτές δεν εξαλείφεται κάθε κίνδυνος. Αρχικά, είναι δυνατή η επαναταυτοποίηση των συμμετεχόντων μέσα από τη χρήση βιολογικών δεδομένων όπως είναι η ηλικία, το φύλο, η φυλή κλπ. Μάλιστα, δεδομένου ότι συχνά υπάρχει οικονομικό κίνητρο για επαναταυτοποίηση των συμμετεχόντων (π.χ. στοχευμένη προώθηση προϊόντων) (110), ο κίνδυνος αυτός πρέπει να λαμβάνεται σοβαρά υπόψη. Οι άνθρωποι που επεξεργάζονται τα δεδομένα καθώς και αυτοί που διαμορφώνουν τις πολιτικές γύρω από αυτά, πρέπει να κάνουν κάθε δυνατή προσπάθεια για να εξασφαλίσουν πως η συλλογή δεδομένων πραγματικού κόσμου, η αποθήκευσή, η κοινοποίηση και η ανάλυσή τους ακολουθούν όλους τις ισχύουσες αρχές. Επίσης, τεχνολογίες και τεχνικές ανάλυσης δεδομένων που προστατεύουν την ιδιωτικότητα πρέπει να χρησιμοποιούνται όπου είναι διαθέσιμες, όπως η διαφορική ιδιωτικότητα (111) και η συνενωμένη μάθηση (112-113). Οι ερευνητές και οι ρυθμιστές θα μπορούσαν να ενσωματώσουν αυτές τις έννοιες στις διαδικασίες συλλογής και ανάλυσης δεδομένων πραγματικού κόσμου καθώς και στην ανακοίνωση των αποτελεσμάτων που προκύπτουν από αυτά ενώ παράλληλα απαιτείται η ολοκλήρωση του νομικού πλαισίου που θα ρυθμίζει τις μελέτες στο πεδίο.

**Διαφορετικότητα, Ισότητα, Αλγοριθμική Δικαιοσύνη και Διαφάνεια:** Οι έννοιες αυτές συνιστούν ένα ακόμα σημαντικό ηθικό ζήτημα γύρω από τα δεδομένα πραγματικού κόσμου. Τα δεδομένα αυτά μπορεί να περιέχουν πληροφορίες από διάφορες δημογραφικές ομάδες οι οποίες μπορούν να αξιοποιηθούν για τη δημιουργία ιατρικών ενδείξεων με βελτιωμένη ικανότητα γενίκευσης συγκριτικά με δεδομένα από τυχαίοποιημένες κλινικές δοκιμές. Από την άλλη, κάποιες πηγές δεδομένων πραγματικού κόσμου μπορούν να είναι επιρρεπείς σε συστηματικά σφάλματα επιλογής και μη ισορροπημένες ως προς κάποιον συγκεκριμένο πληθυσμό με αποτέλεσμα τελικά να έχουν μικρότερη ικανότητα γενίκευσης των αποτελεσμάτων τους (πχ. ορισμένες φορητές

συσκευές ή η πρόσβαση σε συγκεκριμένες εγκαταστάσεις μπορεί να είναι περιορισμένες σε λίγες κοινωνικές ομάδες). Χρειάζεται σημαντική προσπάθεια για να αποκτηθεί πρόσβαση σε δεδομένα από υποαντιπροσωπευόμενους πληθυσμούς. Το θέμα αυτό συνδέεται άμεσα και με την αλγοριθμική δικαιοσύνη η οποία στοχεύει στην κατανόηση και την πρόληψη συστηματικών σφαλμάτων σε μοντέλα μηχανικής μάθησης. Η αλγοριθμική δικαιοσύνη είναι ένα θέμα με αυξημένη ενδιαφέρον στη διεθνή βιβλιογραφία (114-118). Ανακριβή και παραπλανητικά συμπεράσματα μπορούν να προκύψουν από εκπαιδευμένα μοντέλα που συστηματικά αγνοούν μία συγκεκριμένη ομάδα.(π.χ. ένας εκπαιδευμένος αλγόριθμος μπορεί να είναι λιγότερο πιθανό να εντοπίσει έναν καρκίνο σε έναν μαυρό ασθενή συγκριτικά με έναν λευκό). Η διαφάνεια εξασφαλίζει ότι οι πάροχοι των δεδομένων γνωρίζουν πώς χρησιμοποιούνται τα δεδομένα τους και για ποιους σκοπούς (119-122). Η διαφάνεια είναι κρίσιμη για την επεξεργασία δεδομένων πραγματικού κόσμου για το χτίσιμο σχέσεων εμπιστοσύνης μεταξύ των ατόμων από τα οποία προέρχονται τα δεδομένα, τους ερευνητές που τα αναλύουν καθώς και τους ρυθμιστές του πλαισίου.

Τα παραπάνω ζητήματα, παρόλο που μπορούν να αντιμετωπιστούν μεμονωμένα, είναι συχνά αδύνατο να αντιμετωπιστούν ταυτόχρονα δεδομένου ότι συχνά οι αξίες που πρέπει να διαφυλαχθούν είναι αλληλοσυγκρουόμενες. Για το λόγο αυτό είναι ζωτικής σημασίας το χτίσιμο σχέσεων εμπιστοσύνης μεταξύ ερευνητών και συμμετεχόντων, στην οποία θα βασίζεται η επίλυση κάθε σύνθετου θέματος που θα προκύπτει (123). Παράλληλα, είναι σημαντικό να διατηρείται ένας διάυλος επικοινωνίας μεταξύ όλων των εμπλεκομένων για την ταχεία ενημέρωση των παρόχων των δεδομένων και την ανανέωση της συγκατάθεσης όποτε αυτό είναι απαραίτητο (110).

## **Συμπεράσματα**

Τα δεδομένα πραγματικού κόσμου παρέχουν μία πολύτιμη και πλούσια πηγή δεδομένων που ξεφεύγει από τους περιορισμούς των παραδοσιακών επιδημιολογικών μελετών, κλινικών δοκιμών και εργαστηριακών πειραμάτων, με σημαντικό χαμηλότερο κόστος για

την άντλησή τους. Αν χρησιμοποιηθούν και αναλυθούν σωστά έχουν την προοπτική για να αποτελέσουν τη βάση για δημιουργία έγκυρων και ορθών ιατρικών ενδείξεων εξοικονομώντας κόστη και χρόνο, συγκριτικά με τις τυχαιοποιημένες κλινικές μελέτες, και να ενισχύσουν την αποτελεσματικότητα των ιατρικών μελετών. Οι διαδικασίες που βελτιώνουν την ποιότητα των δεδομένων πραγματικού κόσμου και ξεπερνούν τους εγγενείς περιορισμούς τους συνέχισουν συνεχώς να αναπτύσσονται και να βελτιώνονται. Με τον ενθουσιασμό, τη δέσμευση, και την επένδυση στα δεδομένα πραγματικού κόσμου από όλους τους εμπλεκόμενους ελπίζουμε ότι σύντομα θα ξεκλειδώσουμε τις πλήρεις δυνατότητές τους.

1. US Food, Drug Administration, et al. Real-world evidence, 05/20/2022. URL <https://www.fda.gov/science-research/science-and-research-special-topics/real-world-evidence>. accessed on Sep 1, 2022.
2. Wikipedia. Real world data, 2/7/2022. URL [https://en.wikipedia.org/wiki/Real\\_world\\_data](https://en.wikipedia.org/wiki/Real_world_data). accessed on March 19, 2022.
3. Annabel A Powell, Linda Power, Samantha Westrop, Kelsey McOwat, Helen Campbell, Ruth Simmons, Mary E Ramsay, Kevin Brown, Shamez N Ladhani, and Gayatri Amirthalingam. Real-world data shows increased reactogenicity in adults after heterologous compared to homologous prime-boost covid-19 vaccination, march- june 2021, england. *Eurosurveillance*, 26(28):2100634, 2021.
4. Paul R Hunter and Julii Suzanne Brainard. Estimating the effectiveness of the pfizer covid-19 bnt162b2 vaccine after a single dose. a reanalysis of a study of 'real-world' vaccination outcomes from israel. *Medrxiv*, 2021.
5. David A Henry, Mark A Jones, Paulina Stehlik, and Paul P Glasziou. Effectiveness of covid-19 vaccines: findings from real world studies. *The Medical Journal of Australia*, 215(4):149, 2021.
6. Josh A Firth, Joel Hellewell, Petra Klepac, Stephen Kissler, Adam J Kucharski, and Lewis G Spurgin. Using a real-world network to model localized covid-19 control strategies. *Nature medicine*, 26(10):1616–1622, 2020.
7. Allison Shapiro, Nicole Marinsek, Ieuan Clay, Benjamin Bradshaw, Ernesto Ramirez, Jae Min, Andrew Trister, Yuedong Wang, Tim Althoff, and Luca Foschini. Characterizing covid-19 and influenza illnesses in the real world via person-generated health data. *Patterns*, 2(1):100188, 2021.
8. Kira F Ahrens, Rebecca J Neumann, Bianca Kollmann, Michael M Plichta, Klaus Lieb, Oliver Tüscher, and Andreas Reif. Differential impact of covid-related lockdown on mental health in germany. *World Psychiatry*, 20(1):140, 2021.
9. Mauricio A Hernandez and Salvatore J Stolfo. Real-world data is dirty: Data cleansing and the merge/purge problem. *Data mining and knowledge discovery*, 2(1):9–37, 1998.
10. Jacqueline Corrigan-Curay, Leonard Sacks, and Janet Woodcock. Real-world evidence and real-world data for evaluating drug safety and effectiveness. *Jama*, 320(9):867–868, 2018.10
11. Amr Makady, Anthonius de Boer, Hans Hillege, Olaf Klungel, Wim Goettsch, et al. What is real-world data? a review of definitions based on literature and stakeholder interviews. *Value in Health*, 20(7):858–865, 2017.
12. Jessica M Franklin and Sebastian Schneeweiss. When and how can real world data analyses substitute for randomized controlled trials? *Clinical Pharmacology & Therapeutics*, 102(6):924–933, 2017.
13. Rebecca A Miksad and Amy P Abernethy. Harnessing the power of real-world evidence (rwe): a checklist to ensure regulatory-grade data quality. *Clinical Pharmacology Therapeutics*, 103(2):202–205, 2018.



14. Melissa D Curtis, Sandra D Griffith, Melisa Tucker, Michael D Taylor, William B Capra, Gillis Carrigan, Ben Holzman, Aracelis Z Torres, Paul You, Brandon Arneri, et al. Development and validation of a high-quality composite real-world mortality endpoint. *Health services research*, 53(6):4460–4476, 2018.
15. Christopher M Booth, Safiya Karim, and William J Mackillop. Real-world data: towards achieving the achievable in cancer care. *Nature Reviews Clinical Oncology*, 16(5):312–325, 2019.
16. Wencheng Sun, Zhiping Cai, Yangyang Li, Fang Liu, Shengqun Fang, and Guoyan Wang. Data processing and text mining technologies on electronic medical records: a review. *Journal of healthcare engineering*, 2018, 2018.
17. Jionglin Wu, Jason Roy, and Walter F Stewart. Prediction modeling using ehr data: challenges, strategies, and a comparison of machine learning approaches. *Medical care*, pages S106–S113, 2010.
18. Taxiarchis Botsis, Gunnar Hartvigsen, Fei Chen, and Chunhua Weng. Secondary use of ehr: data quality issues and informatics opportunities. *Summit on Translational Bioinformatics*, 2010:1, 2010.
19. Emily Kawaler, Alexander Cobian, Peggy Peissig, Deanna Cross, Steve Yale, and Mark Craven. Learning to predict post-hospitalization vte risk from ehr data. In *AMIA annual symposium proceedings*, volume 2012, page 436. American Medical Informatics Association, 2012.
20. Benjamin Shickel, Patrick James Tighe, Azra Bihorac, and Parisa Rashidi. Deep ehr: a survey of recent advances in deep learning techniques for electronic health record (ehr) analysis. *IEEE journal of biomedical and health informatics*, 22(5):1589–1604, 2017.11
21. Canelle Poirier, Yulin Hswen, Guillaume Bouzillé, Marc Cuggia, Audrey Lavenue, John S Brownstein, Thomas Brewer, and Mauricio Santillana. Influenza forecasting for french regions combining ehr, web and climatic data sources with a machine learning ensemble approach. *PloS one*, 16(5):e0250890, 2021.
22. Tao Zheng, Wei Xie, Liling Xu, Xiaoying He, Ya Zhang, Mingrong You, Gong Yang, and You Chen. A machine learning-based framework to identify type 2 diabetes through electronic health records. *International journal of medical informatics*, 97:120–127, 2017.
23. Rimma Pivovarov, Adler J Perotte, Edouard Grave, John Angiolillo, Chris H Wiggins, and Noémie Elhadad. Learning probabilistic phenotypes from heterogeneous ehr data. *Journal of biomedical informatics*, 58:156–165, 2015.
24. Di Zhao and Chunhua Weng. Combining pubmed knowledge and ehr data to develop a weighted bayesian network for pancreatic cancer prediction. *Journal of biomedical informatics*, 44(5):859–868, 2011.
25. Yogasudha Veturi, Anastasia Lucas, Yuki Bradford, Daniel Hui, Scott Dudek, Elizabeth Theusch, Anurag Verma, Jason E Miller, Iftikhar Kullo, Hakon Hakonarson, et al. A unified framework identifies new links between plasma lipids and diseases from electronic medical records across large-scale cohorts. *Nature genetics*, 53(7):972–981, 2021.

26. Bum Chul Kwon, Min-Je Choi, Joanne Taery Kim, Edward Choi, Young Bin Kim, Soonwook Kwon, Jimeng Sun, and Jaegul Choo. Retainvis: Visual analytics with interpretable and interactive recurrent neural networks on electronic medical records. *IEEE transactions on visualization and computer graphics*, 25(1):299–309, 2018.
27. Elham Mahmoudi, Neil Kamdar, Noa Kim, Gabriella Gonzales, Karandeep Singh, and Akbar K Waljee. Use of electronic medical records in development and validation of risk prediction models of hospital readmission: systematic review. *bmj*, 369, 2020.
28. Rishi J Desai, Shirley V Wang, Muthiah Vaduganathan, Thomas Evers, and Sebastian Schneeweiss. Comparison of machine learning methods with traditional models for use of administrative claims with electronic medical records to predict heart failure outcomes. *JAMA network open*, 3(1):e1918962–e1918962, 2020.
29. Li Huang, Andrew L Shea, Huining Qian, Aditya Masurkar, Hao Deng, and Dianbo Liu. Patient clustering improves efficiency of federated machine learning to predict mortality and hospital stay time using distributed electronic medical records. *Journal of biomedical informatics*, 99:103291, 2019.
30. Victoria L Bartlett, Sanket S Dhruva, Nilay D Shah, Patrick Ryan, and Joseph S Ross. Feasibility of using real-world data to replicate clinical trial evidence. *JAMA network open*, 2(10):e1912869–e1912869, 2019.
31. Nancy A Dreyer and Sarah Garner. Registries for robust evidence. *Jama*, 302(7): 790–791, 2009. 12
32. Stefan Larsson, Peter Lawyer, Göran Garellick, Bertil Lindahl, and Mats Lundström. Use of 13 disease registries in 5 countries demonstrates the potential to use outcome data to improve health care’s value. *Health Affairs*, 31(1):220–227, 2012.
33. Patricia McGettigan, Carla Alonso Olmo, Kelly Plueschke, Mireia Castillon, Daniel Nogueras Zondag, Priya Bahri, Xavier Kurz, and Peter GM Mol. Patient registries: an underused resource for medicines evaluation. *Drug safety*, 42(11):1343–1351, 2019.
34. Peter M Izmirly, Hilary Parton, Lu Wang, W Joseph McCune, S Sam Lim, Cristina Drenkard, Elizabeth D Ferucci, Maria Dall’Era, Caroline Gordon, Charles G Helmick, et al. Prevalence of systemic lupus erythematosus in the united states: estimates from a meta-analysis of the centers for disease control and prevention national lupus registries. *Arthritis & Rheumatology*, 73(6):991–996, 2021.
35. Marijke C Jansen-Van Der Weide, Charlotte MW Gaasterland, Kit CB Roes, Caridad Pontes, Roser Vives, Arantxa Sancho, Stavros Nikolakopoulos, Eric Vermeulen, and Johanna H Van Der Lee. Rare disease registries: potential applications towards impact on development of new drug treatments. *Orphanet Journal of Rare Diseases*, 13(1): 1–11, 2018.
36. Paul Lacaze, Nicole Millis, Megan Fookes, Yvonne Zurynski, Adam Jaffe, Matthew Bellgard, Ingrid Winship, John McNeil, and Alan H Bittles. Rare disease registries: a call to action. *Internal medicine journal*, 47(9):1075–1079, 2017.

37. Bonnie L Svarstad, Theresa I Shireman, and JK Sweeney. Using drug claims data to assess the relationship of medication adherence with hospitalization and costs. *Psychiatric Services*, 52(6):805–811, 2001.
38. Hector S Izurieta, Xiyuan Wu, Yun Lu, Yoganand Chillarige, Michael Wernecke, Arnstein Lindaas, Douglas Pratt, Thomas E MaCurdy, Steve Chu, Jeffrey Kelman, et al. Zostavax vaccine effectiveness among us elderly using real-world evidence: Addressing unmeasured confounders by using multiple imputation after linking beneficiary surveys with medicare claims. *Pharmacoepidemiology and Drug Safety*, 28(7):993–1001, 2019.
39. Alina M Allen, Holly K Van Houten, Lindsey R Sangaralingham, Jayant A Talwalkar, and Rozalina G McCoy. Healthcare cost and utilization in nonalcoholic fatty liver disease: real-world data from a large us claims database. *Hepatology*, 68(6):2230–2238, 2018.
40. Rosarin Sruamsiri, Kosuke Iwasaki, Wentao Tang, and Jörg Mahlich. Persistence rates and medical costs of biological therapies for psoriasis treatment in japan: a real-world data study using a claims database. *BMC dermatology*, 18(1):1–11, 2018.
41. Tiffany P Quock, Tingjian Yan, Eunice Chang, Spencer Guthrie, and Michael S Broder. Epidemiology of al amyloidosis: a real-world study using us claims data. *Blood advances*, 2(10):1046–1053, 2018. 13
42. Matthew Herland, Richard A Bauder, and Taghi M Khoshgoftaar. Medical providerspecialty predictions for the detection of anomalous medicare insurance claims. In 2017 IEEE international conference on information reuse and integration (IRI), pages 579–588. IEEE, 2017.
43. Kenji Momo, Haruna Kobayashi, Yuuka Sugiura, Takeo Yasu, Masayoshi Koinuma, and Sei-ichiro Kuroda. Prevalence of drug–drug interaction in atrial fibrillation patients based on a large claims data. *PLoS One*, 14(12):e0225297, 2019.
44. M Ghiani, U Maywald, T Wilke, and B Heeg. Rw1 bridging the gap between clinical trials and real world data: Evidence on replicability of efficacy results using german claims data. *Value in Health*, 23:S757–S758, 2020.
45. Melih Kirlidog and Cuneyt Asuk. A fraud detection approach with data mining in health insurance. *Procedia-Social and Behavioral Sciences*, 62:989–994, 2012.
46. Jing Li, Kuei-Ying Huang, Jionghua Jin, and Jianjun Shi. A survey on statistical methods for health care fraud detection. *Health care management science*, 11(3):275–287, 2008.
47. Stijn Viaene, Guido Dedene, and Richard A Derrig. Auto claim fraud detection using bayesian learning neural networks. *Expert systems with applications*, 29(3):653–666, 2005.
48. Clifton Phua, Vincent Lee, Kate Smith, and Ross Gayler. A comprehensive survey of data mining-based fraud detection research. *arXiv preprint arXiv:1009.6119*, 2010.
49. Nicolas Roche, Mark Small, Sarah Broomfield, Victoria Higgins, and Ryan Pollard. Real world copd: association of morning symptoms with clinical and patient reported outcomes. *COPD: Journal of Chronic Obstructive Pulmonary Disease*, 10(6):679–686, 2013.
50. M Small, P Anderson, A Vickers, S Kay, and S Fermer. Importance of inhaler-device satisfaction in asthma treatment: real-world observations of physician-observed

- compliance and clinical/patient-reported outcomes. *Advances in therapy*, 28(3):202–212,2011.
51. Jordan E Pinsker, Lars M Müller, Alexandra Constantin, Scott Leas, Michelle Manning, Molly McElwee Malloy, Harsimran Singh, and Steph Habif. Real-world patient-reported outcomes and glycemic results with initiation of control-iq technology. *Diabetes technology & therapeutics*, 23(2):120–127, 2021.
  52. Zahi Touma, Ben Hoskin, Christian Atkinson, David Bell, Olivia Massey, Jennifer H Lofland, Pamela Berry, Chetan S Karyekar, and Karen H Costenbader. Systemic lupus erythematosus symptom clusters and their association with patient-reported outcomes and treatment: Analysis of real-world data. *Arthritis Care & Research*, 2020. 14
  53. Gonzalo J Martinez, Stephen M Mattingly, Shayan Mirjafari, Subigya K Nepal Andrew T Campbell, Anind K Dey, and Aaron D Striegel. On the quality of real-world wearable data in a longitudinal study of information workers. In *2020 IEEE International Conference on Pervasive Computing and Communications Workshops (PerComWorkshops)*, pages 1–6. IEEE, 2020.
  54. Miguel A Hernán, James M Robins, et al. Per-protocol analyses of pragmatic trials. *N Engl J Med*, 377(14):1391–1398, 2017.
  55. Eleanor J Murray, Sonja A Swanson, and Miguel A Hernán. Guidelines for estimating causal effects in pragmatic randomized trials. *arXiv preprint arXiv:1911.06030*, 2019.
  56. Adrian F Hernandez, Rachael L Fleurence, and Russell L Rothman. The adaptable trial and pcorntet: shining light on a new research paradigm, 2015.
  57. Colin Baigent. Pragmatic trials-need for adaptable design. *New England Journal of Medicine*, 384(21), 2021.
  58. Miguel A Hernán and James M Robins. Using big data to emulate a target trial when a randomized trial is not available. *American journal of epidemiology*, 183(8):758–764, 2016.
  59. George N Ioannou, Emily R Locke, Ann M O’Hare, Amy SB Bohnert, Edward J Boyko, Denise M Hynes, and Kristin Berry. Covid-19 vaccination effectiveness against infection or death in a national us health care system: a target trial emulation study. *Annals of Internal Medicine*, 175(3):352–361, 2022.
  60. Xabier García-Albéniz, John Hsu, and Miguel A Hernán. The value of explicitly emulating a target trial when using real world evidence: an application to colorectal cancer screening. *European journal of epidemiology*, 32(6):495–500, 2017. 15
  61. Yoshinori Takeuchi, Hiraku Kumamaru, Yasuhiro Hagiwara, Hiroki Matsui, Hideo Yasunaga, Hiroaki Miyata, and Yutaka Matsuyama. Sodium-glucose cotransporter-2 inhibitors and the risk of urinary tract infection among diabetic patients in japan: Target trial emulation using a nationwide administrative claims database. *Diabetes, Obesity and Metabolism*, 23(6):1379–1388, 2021.
  62. Emily Y Jen, Qing Xu, Aaron Schetter, Donna Przepiorka, Yuan Li Shen, Donna Roscoe, Rajeshwari Sridhara, Albert Deisseroth, Reena Philip, Ann T Farrell, et al. Fda approval: blinatumomab for patients with b-cell precursor acute lymphoblastic leukemia in

- morphologic remission with minimal residual disease. *Clinical Cancer Research*, 25(2):473–477, 2019.
63. Andrea M Gross. Using real world data to support regulatory approval of drugs in rare diseases: A review of opportunities, limitations & a case example. *Current Problems in Cancer*, 45(4):100769, 2021.
64. Jasmanda Wu, Cunlin Wang, Sengwee Toh, Federica Edith Pisa, and Larry Bauer. Use of real-world evidence in regulatory decisions for rare diseases in the united states—current status and future directions. *Pharmacoepidemiology and Drug Safety*, 29(10):1213–1218, 2020.
65. Robin Z Hayeems, Christine Michaels-Igbokwe, Viji Venkataramanan, Taila Hartley, Meryl Acker, Meredith Gillespie, Wendy J Ungar, Roberto Mendoza-Londona, Francois P Bernier, Kym M Boycott, et al. The complexity of diagnosing rare disease: An organizing framework for outcomes research and health economics based on real-world evidence. *Genetics in Medicine*, 24(3):694–702, 2022.
66. Miguel A Hern´an and James M Robins. *Causal inference*, 2010.
67. Martin Ho, Mark van der Laan, Hana Lee, Jie Chen, Kwan Lee, Yixin Fang, Weili He, Telba Irony, Qi Jiang, Xiwu Lin, et al. The current landscape in biostatistics of real-world data and evidence: Causal inference frameworks for study design and analysis. *Statistics in Biopharmaceutical Research*, pages 1–14, 2021.
68. William H Crown. Real-world evidence, causal inference, and machine learning. *Value in Health*, 22(5):587–592, 2019.
69. Peng Cui, Zheyang Shen, Sheng Li, Liuyi Yao, Yaliang Li, Zhixuan Chu, and Jing Gao. Causal inference meets machine learning. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3527–3528, 2020.
70. Hui Y Xiong, Babak Alipanahi, Leo J Lee, Hannes Bretschneider, Daniele Merico, Ryan KC Yuen, Yimin Hua, Serge Gueroussov, Hamed S Najafabadi, Timothy R Hughes, et al. The human splicing code reveals new insights into the genetic determinants of disease. *Science*, 347(6218), 2015.16
71. Daniel Quang, Yifei Chen, and Xiaohui Xie. Dann: a deep learning approach for annotating the pathogenicity of genetic variants. *Bioinformatics*, 31(5):761–763, 2015.
72. Marios Anthimopoulos, Stergios Christodoulidis, Lukas Ebner, Andreas Christe, and Stavroula Mougiakakou. Lung pattern classification for interstitial lung diseases using a deep convolutional neural network. *IEEE transactions on medical imaging*, 35(5): 1207–1216, 2016.
73. Mark JJP Van Grinsven, Bram van Ginneken, Carel B Hoyng, Thomas Theelen, and Clara I S´anchez. Fast convolutional neural network training using selective data sampling: Application to hemorrhage detection in color fundus images. *IEEE transactionson medical imaging*, 35(5):1273–1284, 2016.

74. Jens Kleesiek, Gregor Urban, Alexander Hubert, Daniel Schwarz, Klaus Maier-Hein, Martin Bendszus, and Armin Biller. Deep mri brain extraction: A 3d convolutional neural network for skull stripping. *NeuroImage*, 129:460–469, 2016.
75. Eli Gibson, Wenqi Li, Carole Sudre, Lucas Fidon, Dzhoshkun I Shakir, Guotai Wang, Zach Eaton-Rosen, Robert Gray, Tom Doel, Yipeng Hu, et al. Niftynet: a deep-learning platform for medical imaging. *Computer methods and programs in biomedicine*, 158:113–122, 2018.
76. Mario Coccia. Deep learning technology for improving cancer care in society: New directions in cancer imaging driven by artificial intelligence. *Technology in Society*, 60: 101198, 2020.
77. Nicholas Bien, Pranav Rajpurkar, Robyn L Ball, Jeremy Irvin, Allison Park, Erik Jones, Michael Bereket, Bhavik N Patel, Kristen W Yeom, Katie Shpanskaya, et al. Deep-learning-assisted diagnosis for knee magnetic resonance imaging: development and retrospective validation of mrnet. *PLoS medicine*, 15(11):e1002699, 2018.
78. Fredrik D Johansson, Jamie E Collins, Vincent Yau, Hongshu Guan, Seoyoung C Kim, Elena Losina, David Sontag, Jacklyn Stratton, Huong Trinh, Jeffrey Greenberg, et al. Predicting response to tocilizumab monotherapy in rheumatoid arthritis: a real-world data analysis using machine learning. *The Journal of Rheumatology*, 48(9):1364–1370, 2021.
79. Daniele Rav`i, Charence Wong, Fani Deligianni, Melissa Berthelot, Javier Andreu Perez, Benny Lo, and Guang-Zhong Yang. Deep learning for health informatics. *IEEE journal of biomedical and health informatics*, 21(1):4–21, 2016.
80. Kenji Suzuki. Overview of deep learning in medical imaging. *Radiological physics and technology*, 10(3):257–273, 2017.
81. Dinggang Shen, Guorong Wu, and Heung-Il Suk. Deep learning in medical image analysis. *Annual review of biomedical engineering*, 19:221–248, 2017. 17
82. Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciampi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I S´anchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.
83. June-Goo Lee, Sanghoon Jun, Young-Won Cho, Hyunna Lee, Guk Bae Kim, Joon Beom Seo, and Namkug Kim. Deep learning in medical imaging: general overview. *Korean journal of radiology*, 18(4):570–584, 2017.
84. Amine Amyar, Romain Modzelewski, Hua Li, and Su Ruan. Multi-task deep learning based ct imaging analysis for covid-19 pneumonia: Classification and segmentation. *Computers in Biology and Medicine*, 126:104037, 2020.
85. Yujin Oh, Sangjoon Park, and Jong Chul Ye. Deep learning covid-19 features oncxr using limited training data sets. *IEEE transactions on medical imaging*, 39(8):2688–2700, 2020.
86. Ezz El-Din Hemdan, Marwa A Shouman, and Mohamed Esmail Karar. Covidx-net:A framework of deep learning classifiers to diagnose covid-19 in x-ray images. *arXiv preprint arXiv:2003.11055*, 2020.

87. Shuo Wang, Yunfei Zha, Weimin Li, Qingxia Wu, Xiaohu Li, Meng Niu, Meiyun Wang, Xiaoming Qiu, Hongjun Li, He Yu, et al. A fully automatic deep learning system for covid-19 diagnostic and prognostic analysis. *European Respiratory Journal*, 56(2), 2020.
88. Ali Abbasian Ardakani, Alireza Rajabzadeh Kanafi, U Rajendra Acharya, Nazanin Khadem, and Afshin Mohammadi. Application of deep learning technique to manage covid-19 in routine clinical practice using ct images: Results of 10 convolutional neural networks. *Computers in biology and medicine*, 121:103795, 2020.
89. US Food and Drug Administration. Proposed regulatory framework for modifications to artificial intelligence/machine learning (ai/ml)-based software as a medical device (samd) - discussion paper and request for feedback, 2019. URL <https://www.fda.gov/files/medical%20devices/published/US-FDA-Artificial-Intelligence-and-Machine-Learning-Discussion-Paper.pdf>. accessed on March 24, 2022.
90. US Food and Drug Administration. Artificial intelligence/machine learning (ai/ml)-based software as a medical device (samd) action plan, 01/2021. URL <https://www.fda.gov/media/145022/download>. accessed on March 24, 2022.
91. International Coalition of Medicines Regulatory Authorities. Informal innovation network horizon scanning assessment report - artificial intelligence, 2021. URL [https://www.icmra.info/drupal/sites/default/files/2021-08/horizon\\_scanning\\_report\\_artificial\\_intelligence.pdf](https://www.icmra.info/drupal/sites/default/files/2021-08/horizon_scanning_report_artificial_intelligence.pdf). accessed on March 24, 2022.
92. European Medicine Agency. Artificial intelligence in medicine regulation, 2021. URL <https://www.ema.europa.eu/en/news/artificial-intelligence-medicine-regulation>. accessed on March 24, 2022. 18
93. Mark J Van der Laan and Sherri Rose. Targeted learning: causal inference for observational and experimental data, volume 4. Springer, 2011.
94. Mark J Van der Laan and Sherri Rose. Targeted learning in data science. Springer, 2018.
95. Mark J van der Laan and Alexander R Luedtke. Targeted learning of the mean outcome under an optimal dynamic treatment rule. *Journal of causal inference*, 3(1):61–95, 2015.
96. P Chakravarti, A Wilson, S Krikov, N Shao, and M van der Laan. Pin68 estimating effects in observational real-world data, from target trials to targeted learning: Example of treating covid-hospitalized patients. *Value in Health*, 24:S118, 2021.
97. Hans-Georg Eichler, Franz Koenig, Peter Arlett, Harald Enzmann, Anthony Humphreys, Frank P´etavy, Brigitte Schwarzer-Daum, Bruno Sepodes, Spiros Vamvakas, and Guido Rasi. Are novel, nonrandomized analytic methods fit for decision making? the need for prospective, controlled, and transparent validation. *Clinical Pharmacology & Therapeutics*, 107(4):773–779, 2020.
98. Supriyo Chakraborty, Richard Tomsett, Ramya Raghavendra, Daniel Harborne, Moustafa Alzantot, Federico Cerutti, Mani Srivastava, Alun Preece, Simon Julier, Raghuvver M Rao, et al. Interpretability of deep learning models: A survey of results. In 2017 IEEE smartworld, ubiquitous intelligence & computing, advanced & trusted computed, scalable computing & communications, cloud & big data computing, Internet of people and smart

- city innovation (smart-world/SCALCOM/UIC/ATC/CBDcom/IOP/SCI), pages 1–6. IEEE, 2017.
99. Quanshi Zhang and Song-Chun Zhu. Visual interpretability for deep learning: a survey. arXiv preprint arXiv:1802.00614, 2018.
  100. Fred Hohman, Haekyu Park, Caleb Robinson, and Duen Horng Polo Chau. Summit: Scaling deep learning interpretability by visualizing activation and attribution summarizations. IEEE transactions on visualization and computer graphics, 26(1):1096–1106, 2019.
  101. Biraja Ghoshal and Allan Tucker. Estimating uncertainty and interpretability in deep learning for coronavirus (covid-19) detection. arXiv preprint arXiv:2003.10769, 2020.
  102. Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. arXiv preprint arXiv:1706.05806, 2017.19
  103. Angel Alfonso Cruz-Roa, John Edison Arevalo Ovalle, Anant Madabhushi, and Fabio Augusto González Osorio. A deep learning architecture for image representation, visual interpretability and automated basal-cell carcinoma cancer detection. In International Conference on Medical Image Computing and Computer-Assisted Intervention, pages 403–410. Springer, 2013.
  104. Lorena A Barba. Terminologies for reproducible research. arXiv preprint arXiv:1802.03311, 2018.
  105. Aaron Stupple, David Singerman, and Leo Anthony Celi. The reproducibility crisis in the age of digital medicine. NPJ digital medicine, 2(1):1–3, 2019.
  106. Rickey E Carter, Zach I Attia, Francisco Lopez-Jimenez, and Paul A Friedman. Pragmatic considerations for fostering reproducible research in artificial intelligence. NPJ digital medicine, 2(1):1–3, 2019.
  107. Chao Liu, Cuiyun Gao, Xin Xia, David Lo, John Grundy, and Xiaohu Yang. On the replicability and reproducibility of deep learning in software engineering. arXiv preprint arXiv:2006.14244, 2020.
  108. Xafis V, Schaefer GO, Labude MK, Brassington I, Ballantyne A, Lim HY, Lipworth W, Lysaght T, Stewart C, Sun S, Laurie GT, Tai ES. An Ethics Framework for Big Data in Health and Research. Asian Bioeth Rev. 2019 Oct 1;11(3):227-254.
  109. Yang KM. Emerging Ethical Issues in Managing Real World Data. Healthc Inform Res. 2019 Oct;25(4):237-238.
  110. Tanner A. *Or bodies, our data: how companies make billions selling our medical records*. Boston (MA): Beacon Press; 2017.
  111. Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In Theory of cryptography conference, pages 265–284. Springer, 2006.
  112. Jakub Konečný, Brendan McMahan, and Daniel Ramage. Federated optimization: Distributed optimization beyond the datacenter. arXiv preprint arXiv:1511.03575, 2015.



113. Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. arXiv preprint arXiv:1610.05492, 2016.20
114. Melissa D McCradden, Shalmali Joshi, Mjaye Mazwi, and James A Anderson. Ethical limitations of algorithmic fairness solutions in health care machine learning. *The Lancet Digital Health*, 2(5):e221–e223, 2020.
115. Shira Mitchell, Eric Potash, Solon Barocas, Alexander D’Amour, and Kristian Lum. Algorithmic fairness: Choices, assumptions, and definitions. *Annual Review of Statistics and Its Application*, 8:141–163, 2021.
116. Vishwali Mhasawade, Yuan Zhao, and Rumi Chunara. Machine learning and algorithmic fairness in public and population health. *Nature Machine Intelligence*, 3(8):659–666, 2021.
117. Pak-Hang Wong. Democratizing algorithmic fairness. *Philosophy & Technology*, 33(2): 225–244, 2020.
118. Jessica K Paulus and David M Kent. Predictably unequal: understanding and addressing concerns that algorithmic clinical prediction may increase health disparities. *NPJ digital medicine*, 3(1):1–8, 2020.
119. Lucinda S Orsini, Marc Berger, William Crown, Gregory Daniel, Hans-Georg Eichler, Wim Goettsch, Jennifer Graff, John Guerino, Pall Jonsson, Nirosha Mahendraratnam Lederer, et al. Improving transparency to build trust in real-world secondary data studies for hypothesis testing—why, what, and how: recommendations and a roadmap from the real-world evidence transparency initiative. *Value in Health*, 23(9):1128–1136, 2020.
120. Elisabetta Patorno, Sebastian Schneeweiss, and Shirley V Wang. Transparency in real-world evidence (rwe) studies to build confidence for decision-making: reporting rwe research in diabetes. *Diabetes, Obesity and Metabolism*, 22:45–59, 2020.
121. Richard White. Building trust in real-world evidence and comparative effectiveness research: the need for transparency, 2017.
122. Elena Rodriguez-Villa and John Torous. Regulating digital health technologies with transparency: the case for dynamic and multi-stakeholder evaluation. *BMC medicine*, 17(1):1–5, 2019.
123. van Staa T-P, Goldacre B, Buchan I, Smeeth L. Big health data: the need to earn public trust. *BMJ*. 2016;354:i3636