

LES DONNÉES DE SANTÉ, UN PATRIMOINE COMMUN QUI DOIT SERVIR À AMÉLIORER LE BIEN-ÊTRE DE TOUS

Exhaustives et précises sur nos conditions de vie, nos soins..., de nombreuses informations sont collectées tous les jours mais restent peu ou pas disponibles pour la recherche, regrette **un collectif de professeurs de médecine**

Dans un contexte de pandémie mondiale, à l'heure où le prix Nobel de chimie d'Emmanuelle Charpentier *[qui mène ses recherches à l'étranger]* questionne l'avenir de la recherche en France, il est urgent de rap- peler que nous disposons d'un immense potentiel sous-exploité : des millions de données administratives, sociales, médi- cales et environnementales.

En tant que médecins, chercheurs, épi- démiologistes et économistes, nous tra- vaillons tous les jours à résoudre des pro- blématiques complexes pour améliorer le bien-être de tous. Nous cherchons à com- prendre pourquoi la France a un écart d'espérance de vie en fonction du milieu social, comme par exemple entre ou- vriers et cadres supérieurs, plus marqué que ses voisins européens ; comment le

Levothyrox, les prothèses mammaires PIP ou le Mediator ont pu rester si long- temps en circulation malgré leur dange- rosité ; pourquoi nos chances de mourir varient selon que l'on fait un AVC dans le Nord ou en Occitanie. Cependant, nous nous heurtons toujours à la même diffi- culté : l'accès limité aux données de santé. Pourtant, en France, des données exhaustives et précises sur nos condi- tions de vie, notre santé, les soins que nous recevons sont collectées tous les jours. Ces données sont nombreuses, souvent de qualité, mais peu ou pas disponibles pour la recherche. Nous ne comptons plus les travaux talents, voire abandonnés, faute d'un accès effectif aux données de santé. Les trois ou quatre ans pouvant être nécessaires à l'obtention d'un simple jeu de données sont to- talement inacceptables.

Dépasser une logique propriétaire

Il est urgent pour notre pays que la re- cherche puisse confronter les données de santé, de recours aux soins, socio-écono- miques et environnementales. De façon très concrète, une actualisation et un croisement rapide des connaissances auraient permis d'identifier plus précocement les facteurs de vulnérabilité vis-à- vis de l'infection du Covid-19, de son évolu- tion et du vécu du confinement. Les dé- bats techniques et polémiques autour de la mobilisation des données de santé au service de la recherche ne doivent pas oc- culter l'enjeu principal : la santé de nos concitoyens et notre capacité à faire face aux grands enjeux présents et futurs qui s'y rapportent. Pour cela, il nous faut

collectivement changer de regard sur ces données. Cesser de les considérer seule- ment comme un instrument de gestion

administrative pour en faire un outil opé- rationnel de recherche. Dépasser une cer- taine logique propriétaire et accepter que leur utilité se révèle dans le partage. Nous devons les voir comme le patrimoine commun des citoyens et ne les utiliser que pour servir la société et la population. Toutes les institutions de notre pays (hôpitaux, mairies, Assurance-mala- die...) génèrent et gèrent des données sous des formes variées, avec des contrô- les de qualité variables. Les organismes producteurs ne communiquent pas, et parfois ne désirent pas collaborer ou par- tager les données avec des tiers, les ren- dant difficilement accessibles. Tout se passe comme s'ils en étaient eux-mêmes propriétaires alors qu'elles appartiennent de fait aux citoyens.

Les raisons invoquées sont parfois légi- times, à titre d'exemple, ils souhaitent être crédités pour leur travail. A contra- rio, il n'est pas rare que d'autres raisons soient avancées, telles que la question de la protection des données ou de leur sé- curité, dans un pays où la réglementation est l'une des plus strictes en la ma- tière, comme un prétexte à réserver ces données à leurs seuls travaux. Veiller à valoriser leurs travaux est indispensable, accepter que les données leur soient ré- servées est impensable. Ces données sont un patrimoine commun devant ser- vir à améliorer la santé de tous. Il est in- dispensable de fournir un cadre éthique, technique, humain et financier à même d'assurer cette mutualisation au service de la recherche. Ce cadre doit répondre à plusieurs impératifs : respecter la vie pri- vée des citoyens ; respecter le travail de tous ceux qui participent à cette œuvre commune ; réserver cette mutualisation à la poursuite de l'intérêt général.

De façon très concrète, nous deman- dons une gouvernance claire, qui aura pour objectifs : rassembler les données

de soins, sociales et environnementales ; faciliter leur accès et leur exploitation ra- pides ; créer un modèle de valorisation pour les chercheurs ou organismes qui collectent et partagent les données ; flé- cher les investissements vers la constitu- tion et la pérennisation de notre patri- moine commun de données, en passant par le renforcement des équipes concer- nées en conditionnant l'attribution des financements publics à leur mise à dispo- sition facile et rapide pour la collectivité.

Permettre un système d'excellence

Le Secur de la santé constitue l'occasion d'un geste fort pour amorcer ce mouve- ment en finançant le développement des entrepôts de données hospitaliers. De même, l'Agence nationale de la re- cherche pourrait prioriser les appels d'of- fres ayant pour objectifs la mutualisa- tion et la valorisation scientifique des données. Enfin, les initiatives d'excel- lence et les trois instituts nationaux d'in- telligence artificielle impliqués en santé sont particulièrement bien placés pour porter des initiatives dans ce domaine.

Ces demandes sont réalistes, nous avons toutes les qualités pour permettre un grand saut qualitatif de nos politiques de santé : un système de santé d'excel- lence, une protection sociale forte, un cadre éthique consensuel et une recher- che de grande qualité. Collectivement né- cessaire pour anticiper et suivre les crises sanitaires, assurer la sécurité des traite- ments et lutter contre les inégalités sociales ; individuellement utile pour prévenir les maladies, permettre des dia- gnostics précoces et une prise en charge adaptée tout au long de la vie, l'acces- sibilité des données de santé à la recher- che n'est pas une question triviale. C'est un déterminant de l'accès aux soins pour la population et de la pertinence de la décision publique, de souveraineté politi- que. C'est une condition du libre exercice du débat démocratique. ■

Karine Chevreul, direc- trice de l'unité Inserm 1123,

professeure de santé publi- que et économie de la santé, université de Paris,

Hôtel-Dieu, AP-HP ; **Cyrille**

Delplâtre, directeur de re- cherche Inserm, épidémi- ologiste, université de Tou- louse-III ; **Paul Dougron**, directeur de recherche, éco- nomiste de la santé à l'Ins- titut de recherche et docu- mentation en économie de la santé (Ileds) ; **Martine**

Gillard, professeure en car- diologie, université de

Brest, CHU de Brest ; **Mi- chelle Kelly-Irving**, chir-

gienne de recherche Inserm, épidémiologiste sociale,

université de Toulouse-III ;

Bertrand Lutkacs, chirur- gien urologue, hôpital Te- non, AP-HP ; **Alexandre**

Mebazza, directeur de l'unité Inserm 942, profes-

seur en anesthésie-réani- mation, université de Paris,

hôpital Lariboisière, AP-HP ;

Jean-Louis Pepin, direc- teur de l'unité Inserm UMR

1042, professeur en physi- ologie clinique, université et

CHU Grenoble-Alpes.

Marcel Goldberg et Marie Zins La plate-forme « Health Data Hub » pose des questions de sécurité majeures

Avec cette infrastructure informatique, qui permet d'héberger et de centraliser les informations de santé des Français à des fins de recherche médicale, il suffit de croiser quelques données simples pour identifier une personne, s'inquiètent les deux épidémiologistes

Le gouvernement a lancé un très ambitieux projet de « Health Data Hub » (HDH) visant à réunir l'ensemble des données disponibles sur la santé des Français, pour dévelop- per l'intelligence artificielle (IA) en santé. En effet, la situation ac- tuelle est largement insatisfai- sante en raison de la dispersion en de multiples systèmes d'informa- tion gérés sans coordination par de nombreux acteurs : hôpi- taux, Sécurité sociale, organis- mes de recherche, universités, re- gistres et enquêtes épidémiologi- ques, cohortes... On ne peut que souscrire aux objectifs de partage de données et de développement de l'IA en santé et se féliciter de la volonté politique de fournir des moyens conséquents.

Mais si l'intention est louable, réunir toutes les données dans une infrastructure informatique unique est extrêmement dange- reux et largement inutile. Le fait de confier sa gestion à Microsoft a suscité de nombreux débats, mais on n'a pratiquement pas évoqué les très graves problèmes que pose le dispositif prévu, même s'il était géré sur une in- frastructure nationale.

En effet, le HDH entend centra- liser toute donnée collectée dans le cadre d'un acte remboursé par l'Assurance-maladie dans les hô- pitaux, en médecine de ville, mé- decine du travail, pharmacies, services de protection maternelle et infantile, dépistage, enquêtes de santé... La centralisation des données concernant les aspects les plus intimes de la vie des 67 millions d'assurés sociaux chez un hébergeur unique pose des questions majeures de sécu- rité, car il suffit de croiser quel- ques données simples pour iden- tifier une personne, avec des con- séquences potentiellement très lourdes. La centralisation des données dans une seule infras- tructure informatique peut per- mettre des mesures de sécurité accrues, mais les rend plus expo- sées aux attaques venant de l'ex- térieur comme de l'intérieur, avec des impacts plus grands en cas de rupture de confidentialité.

Faire courir un tel danger aux personnes ne peut se justifier que si cela est indispensable. Or ce n'est pas le cas : non seulement un système centralisé est dange- reux, mais il est largement inutile pour deux raisons essentielles.

La première tient à la qualité des différentes bases de données con- cernées. Construites dans des buts, des circonstances et avec des méthodes qui, pour la plupart, n'ont rien à voir entre elles, leur qualité et leur validité sont extrê- mement variables : « big data » n'est pas synonyme de « good data ». Or les algorithmes d'intelli- gence artificielle ont besoin de données valides. Avant d'utiliser une base de données, un examen minutieux de ses caractéristiques et de sa qualité, impliquant ceux qui l'ont construite, est indispen- sable, sans quoi son intégration dans le HDH est inutile.

Aberration scientifique

La seconde raison est l'hétérogé- nité de ces bases de données. Il ne suffit pas de regrouper des données de droite et de gauche pour les « faire parler ». Encore faut-il que les données soient in- teropérables, c'est-à-dire homo- gènes sur le plan sémantique. Par exemple, si on s'intéresse à l'in- suffisance cardiaque, on peut trouver des données dans diver- ses sources : dossier de service de cardiologie, diagnostic de généra- liste ou de cardiologue en ville,

déclaration d'un sujet dans une enquête, réseaux sociaux... Mais, selon la source, ce terme n'a pas la même signification ni la même validité. Il faut connaître le con- texte et les méthodes du recueil des données, la population dont elles sont issues, etc., le cas échéant le type d'appareil utilisé car on rencontre, par exemple, des électrocardiogrammes ou des images IRM provenant d'ap- pareils différents. Et dans de nombreux cas, cette harmonisa- tion s'avère impossible. C'est pourquoi les algorithmes d'IA sont le plus souvent développés sur une base de données unique.

Il arrive cependant que plu- sieurs bases de données puissent être rassemblées. Il faut alors les harmoniser. Mais ceci n'a de sens que pour des objectifs spéci- fiques de recherche et ne peut donc être réalisé qu'au cas par cas ; et implique un travail de comparaison et de définition des données, qui ne peut être réalisé que par les responsables des don- nées concernées, qui disposent de l'expertise et de la connais- sance approfondie des données, des conditions de leur recueil, des modalités de validation...

Imaginer qu'il sera possible de développer des algorithmes d'IA à partir des données extrême- ment hétérogènes uniquement parce qu'elles sont stockées dans un système informatique centra- lisé est donc une aberration scientifique et technique.

Analyse « distribuée »

Et même si toutes ces difficultés sont résolues et qu'on dispose de plusieurs bases de données véri- tablement interopérables, il n'est pas indispensable de les réunir dans la même infrastructure in- formatique. Il existe des métho- des d'analyse « distribuée » où des données gérées dans des sys- tèmes informatiques différents sont exploitées en commun : ces méthodes sont largement utili- sées dans les cas où, pour des ra- sons de sécurité ou des raisons lé- gales, les données ne doivent pas être transférées hors de leur pro- pre environnement.

Développer le HDH en réunis- sant les données de 67 millions de Français dans une infrastructure informatique unique est donc une erreur fondamentale qui fait inutilement courir de graves dan- gers. Si les objectifs de partage de

données et de développement de l'IA sont pleinement justifiés, plu- tôt que d'empêcher aveuglément des bases de données hétérocli- tes, le HDH devrait se concentrer sur des activités réellement utiles, comme par exemple une carto- graphie analytique des bases de données disponibles. Leur mise en réseau la promotion de règles harmonisées de partage de don- nées, etc. En gardant à l'esprit que les véritables difficultés se situent à la source même des données de santé, comme la crise sanitaire du Covid-19 l'a cruellement mis en évidence : absence de données provenant des Ehpad, insuffi- sance du nombre de spécialistes du codage des causes de décès, pour ne citer que les manques les plus voyants. ■

Marcel Goldberg est profes- seur émérite d'épidémiologie et de santé publique à l'univer- sité Paris-Descartes ; **Marie**

Zins est médecin épidémiolo- giste, enseignante-chercheuse à l'université Paris-Descartes

