

Les limites des modèles comportementaux du big data

Dans sa chronique mensuelle « Recherches », l'économiste Paul Seabright explique pourquoi il ne faut pas surestimer la qualité de la compréhension qu'on peut tirer des « données massives ».

LE MONDE ECONOMIE | 16.03.2017 à 11h06 • Mis à jour le 16.03.2017 à 15h06 | Par Paul Seabright (Institut d'études avancées de Toulouse)



« Les entreprises comme Google, Apple ou Facebook disposent de millions de variables décrivant le comportement de millions de personnes. Les techniques d'apprentissage machine ("machine learning") peuvent y cerner des tendances qui échapperaient à un regard purement humain ». KamiPhuc/Flickr/CC BY 2.0

Le débat sur les avantages et les inconvénients des big data (« données massives ») tend à opposer les bénéfices d'une meilleure compréhension des comportements humains aux dangers d'abus concernant la vie privée. Un article de Susan Athey, ancienne économiste en chef de Microsoft, montre qu'il ne faut pas non plus surestimer la qualité de la compréhension qu'on peut en tirer (« Beyond Prediction : Using Big Data for Policy Problems », *Science* n° 6324, 3 février 2017).

Jusqu'à récemment, les analyses statistiques des comportements humains devaient choisir entre deux types d'informations. Les enquêtes permettent de poser beaucoup de questions à relativement peu de gens, avec le risque que les personnes interrogées soient peu représentatives de l'ensemble de la population. Avec un échantillon plus large, les recensements permettent de s'adresser à beaucoup de gens, voire à des populations entières, mais en leur posant peu de questions, ce qui limite l'analyse à une modélisation simple.

Mise en garde

Désormais, les entreprises comme Google, Apple ou Facebook disposent de millions de variables décrivant le comportement de millions de personnes. Les techniques d'apprentissage machine (« *machine learning* ») peuvent y cerner des tendances qui échapperaient à un regard purement humain. Mais Susan Athey nous met en garde contre un optimisme facile quant à la sophistication des modèles comportementaux qui en découlent.

Observer des comportements, les cerner à l'aide de l'algorithme le plus sophistiqué, ne nous aide pas à savoir si ces comportements restent inchangés lorsque nous essayons d'intervenir pour améliorer la situation. Or, quasiment toutes les applications des big data concernent une intervention potentielle, que ce soit une politique publique, la politique commerciale d'une entreprise ou le choix d'un hôpital entre différents traitements.

Son article cite de nombreux cas où les tendances observées par les méthodes du big data ne suffisent pas pour prédire l'impact d'une intervention. La société eBay avait cru calculer par ces méthodes que son retour sur investissement en publicité en ligne était de 1 400 % à cause d'une forte corrélation entre les achats et les investissements publicitaires. Après une vérification expérimentale, il a été constaté que le vrai retour était de... - 63 %, car la plupart des achats auraient été faits sans les annonces !

Des risques scientifiques

L'apprentissage machine est souvent utilisé par les entreprises privées pour prédire les profils de clients les plus susceptibles de quitter la firme pour un concurrent. Ces prédictions sont utilisées pour allouer le service après-vente en priorité aux clients de fidélité faible. Mais ces interventions sont souvent décevantes : être susceptible de partir vers un concurrent ne rend pas forcément le client sensible aux efforts de la firme pour le garder.

Un exemple des risques scientifiques des analyses big data apparaît dans un article d'un autre économiste en chef d'une grande entreprise – en l'occurrence Hal Varian, de Google (« **Big Data : New tricks for Econometrics** », *Journal of Economic Perspectives* (<https://www.aeaweb.org/articles?id=10.1257/jep.28.2.3>), n°28/2, printemps 2014). On constate depuis vingt ans qu'être noir aux Etats-Unis est associé à une probabilité plus faible de se voir accorder un prêt immobilier. Une analyse big data effectuée par Varian montre que lorsqu'on prend en compte le fait d'avoir pu trouver ou non une assurance (condition nécessaire pour un prêt), la différence raciale ne joue plus aucun rôle.

Peut-on en conclure que les différences raciales ne sont pas importantes pour accéder aux prêts immobiliers ? Pas du tout ! Comme le reconnaît Varian, trouver une assurance pourrait être plus difficile pour les Américains noirs que pour les autres – c'est peut-être même à travers l'allocation des assurances que la discrimination raciale aurait son impact principal sur l'accès aux prêts. Les big data permettent de prédire qui recevra un prêt, mais en expliquer les causes reste un défi autrement plus complexe.