

## Algorithmes 2/3

<p>Soulever le capot</p> <p>Sortir de la phobie ou de la phobie</p> <p>2 parties :</p> <ul style="list-style-type: none"> <li>-les systèmes à apprentissage</li> <li>-les données</li> </ul>	<p>Avertissement : cette séance va traiter d'autres types d'algorithmes pour lesquels, comme pour APB, on va s'interroger sur les impacts qu'ils peuvent avoir sur nous. Mais, comme pour APB, et selon l'intention exprimée en introduction avec le texte de Dominique Cardon, il va falloir « ouvrir le capot ». Dans la première séance j'avais du parler de numérique, d'information, de calculateur, ... cette fois on va parler de corrélation, de réseaux de neurones, d'apprentissage, de big data.... Ce sera évidemment sommaire mais ce me semble indispensable pour comprendre ces automatismes et pouvoir ainsi les analyser et, le cas échéant, les critiquer.</p>
<p>Différentes familles d'algo</p>	<ul style="list-style-type: none"> <li>- Les algorithmes d'appariement (matching) : très répandus : école, emploi, médecine, rencontres, ....(pas de prix= Les questions à se poser sont du même ordre que celles envisagées pour APB.</li> <li>- Les algorithmes de recommandation : de livres, de films, de formation, d'emploi, de vote, etc. Les questions importantes : les données de départ, la bulle ?, Loyauté difficile.</li> <li>- Les algorithmes d'aide à la décision : GPS, justice, médecine, police, Poids donné au passé (ex : liberté conditionnelle) et PB de responsabilité.</li> <li>- Les algorithmes de prédiction : décrochage scolaire, épidémie, anti-terrorisme, accueil des œuvres culturelles,... PB de choix des critères, des données, ....</li> <li>- Les algorithmes générateurs de connaissances : page rank de Google, logiciel de l'affaire Grégory, ...</li> </ul>
<p><b>A croiser avec la technique</b></p> <p>Déterminisme</p> <p>Logique déductive</p> <p>Basé sur des études consultables et critiquables</p>	<p>L'algorithme utilisé par APB repose certes sur la trouvaille des deux mathématiciens (qui leur a valu le prix Nobel) mais l'idée une fois trouvée, il n'y a derrière qu'une programmation informatique tout à fait ordinaire</p> <p>L'algorithme utilisé ici est <b>déterministe</b> : il doit donner une solution. ou, au minimum, plusieurs solutions possibles entre lesquelles choisir. La plupart des algorithmes insérés dans le monde technique le sont aussi. Le problème à résoudre peut être compliqué (aller de Clermont à Barcelone en transport en commun...) mais les éléments du problème existent.</p> <p>Dans ce type d'algorithme, on est dans le domaine de la logique déductive : la réponse est certaine. Les résultats sont déduits des règles et/ou des résultats précédents : si l'élève a telles ou telles caractéristiques alors il est affecté à telle formation ; si le parcours x est plus court ou plus rapide que le parcours y alors choisir x, etc.</p> <p>Les caractéristiques qui vont déterminer le résultat du calcul peuvent être complexes, résultats accumulés de travaux, de connaissances, etc. Si l'expérience du magistrat qui doit décider ou non de la liberté conditionnelle d'un détenu observe que cette personne n'a ni emploi, ni qualification, ni famille, ni logement....son expérience ou les consignes qu'on lui donne peuvent lui faire refuser la mise en liberté. Si une étude des mises en liberté conditionnelles a été effectuée, des consignes ont pu être énoncées avec les explications adéquates. Ces consignes peuvent être incorporées dans un programme informatique ce qui donnera un algorithme d'aide à la décision du magistrat. On verra que cette automatisation provoquée par le passage de l'étude sociologique à l'algorithme pose déjà une série de problèmes (quelle est alors la liberté et la responsabilité du juge ?) mais l'étude</p>

	sociologique existe, on peut la consulter, la critiquer, savoir pourquoi et dans quelle mesure on peut suivre l'algorithme....
<b>Nouveauté radicale</b>  Logique inductive à partir des data	La nouveauté radicale des algorithmes par apprentissage automatique c'est qu'ils vont partir directement des faits pour induire les règles à appliquer. On passe de la logique déductive à une logique inductive. En soi où est le mal ? Nous pratiquons nous-mêmes ce type de logique : lorsqu'un même fait se répète on est tenté de penser qu'il va se reproduire. Par exemple, je sais qu'il arrive fréquemment que telle personne arrive en retard, j'agis en conséquence MAIS je sais que ce n'est pas certain : cette fois, elle peut être à l'heure. Où est le problème, sinon forcément le mal, c'est que l'algorithme va élaborer lui-même ses règles à partir des faits. En restant dans l'exemple de la remise de peine, il <u>va d'abord remplacer l'étude des sociologues</u> .

### Traduction automatique

Pour illustrer cette évolution dans la méthode des algorithmes, prenons un exemple qui ne soit pas trop chargé d'humanité, celui de la traduction automatique d'une langue dans une autre.

1<sup>ère</sup> approche : Dès les années 50, des linguistes et informaticiens, pour élaborer un algorithme de traduction, prennent un dictionnaire, une grammaire dont ils extraient des règles pour construire un modèle (cf. les travaux de Jacobson, Noam Chomsky, ...). Ça donnera SYSTRAN. Mais il reste beaucoup de problèmes (exemple, the « flesh is weak » est traduit « par la viande est molle »).

2<sup>ème</sup> approche : à partir de 80, au lieu de travailler comme un traducteur, la machine, recevant une phrase à traduire, va plonger dans des « corpus » bilingues (dans le cas de l'anglais, provenant du Canada ou de l'Union Européenne) pour trouver la phrase qui est la plus fréquemment trouvée. Aujourd'hui, dans GoogleTrad, « flesh is weak » est traduit par « la chair est faible » : c'est l'exploitation de la « force brute » de l'ordinateur, rapidité d'exploration du fonds de référence + un peu de statistique. (voir différence avec le Gaffiot où le mauvais élève cherche une bonne traduction alors que le programme va choisir la traduction la plus fréquente).

3<sup>ème</sup> approche : en 2014 : application des réseaux de neurones artificiels : on transforme les textes en vecteurs sur lesquels on effectue les traitements (un pas de plus dans la numérisation !). On fait encore lire à la machine des millions de pages dans les deux langues mais le but n'est plus de traduire mais d'apprendre à traduire. Une fois l'apprentissage réalisé, on n'a plus besoin des bases d'apprentissage, la machine traduit selon le « modèle » (« pattern ») qu'elle s'est fait elle-même (ou presque<sup>1</sup>) par son apprentissage.... Modèle qui n'a rien à voir avec celui que recherchaient les linguistes-informaticiens dans les années 50

Et ça fonctionne ! Google, Systran, Facebook, DeepL...les grandes entreprises de traduction automatiques sont passées ou sont en train de passer à cette nouvelle méthode.

(Pour en savoir plus sur [http://www.lemonde.fr/sciences/article/2017/11/27/la-traduction-dopee-par-l-intelligence-artificielle\\_5221041\\_1650684.html#lwRTavLiT7atL4dL.99](http://www.lemonde.fr/sciences/article/2017/11/27/la-traduction-dopee-par-l-intelligence-artificielle_5221041_1650684.html#lwRTavLiT7atL4dL.99) )

<sup>1</sup> En fait, l'apprentissage est « supervisé » par des experts du domaine mis en apprentissage



<p><b>Les réseaux de neurones</b></p> <p>-un bon niveau en maths <del>en biologie</del>          -le procédé s'applique à une grande variété de domaines,          -ça marche mais les créateurs ne savent pas toujours ni pourquoi ni comment ! (Il y a comme cela des domaines des mathématiques où subsistent quelque mystère cf. par exemple, les automates cellulaires de J.Conway)  <a href="https://www.youtube.com/watch?v=S-W0NX97DB0">https://www.youtube.com/watch?v=S-W0NX97DB0</a>   <a href="http://jean-paul.davalan.pagesperso-orange.fr/divers/jeuvie/index.html">http://jean-paul.davalan.pagesperso-orange.fr/divers/jeuvie/index.html</a></p>	<p>Cf NOTES p.3</p> <p>Une présentation par l'un de ses inventeurs, professeur au collège de France, Yann LeCun : <a href="http://www.youtube.com/watch?v=OzZoPVjv8iE">www.youtube.com/watch?v=OzZoPVjv8iE</a></p> <p>Plus technique : <a href="https://www.imagile.fr/creation-dun-petit-reseau-de-neurones-artificiels-netait-complexe/">https://www.imagile.fr/creation-dun-petit-reseau-de-neurones-artificiels-netait-complexe/</a></p>
---	--

**« C'est fascinant de voir que cette technique, qui reste encore opaque et mal comprise, fonctionne aussi bien », constate François Yvon<sup>2</sup> ».**

*On peut avoir une bonne illustration des possibilités et de ces réseaux de neurones artificiels avec le film « Alpha-GO » visible sur NETFLIX qui retrace la préparation du programme qui, en 2017, a battu le champion du monde du jeu de GO, Se-Dol Lee. (Netflix offre un mois d'essai gratuit : ça peut servir à ça !)*

*En 1997, un ordinateur, Deep Blue, battait Kasparov, champion du monde des échecs. En 2017, le programme Alpha Go battait Se-Dol Lee, champion du monde de Go.*

*-IBM en 97, Google en 2017,*

*-Un ordinateur de 1,5T en 97, un programme dans un portable en 2017,*

*-Exploitation de toutes les parties jouées en 97, apprentissage du modèle du jeu de Go en 2017,*

*-étonnement du champion devant les coups inventés par le programme mais l'étonnement tout autant des ingénieurs sidérés des réactions de leur programme.*

*-Lorsque l'on dit qu'un algorithme est incapable de créer, il faut s'entendre sur ce qu'est une création : il ne peut créer à partir de rien mais, ayant appris un « modèle » il est capable de produire des résultats (ou, ici, des coups) auxquels nul n'avait encore pensé.*

---

<sup>2</sup> François Yvon, directeur du Laboratoire d'informatique pour la mécanique et les sciences de l'ingénieur (Limsi-CNRS) à Orsay (Essonne)

## Mais l'apprentissage se fait à partir des données : ce sont elles qui vont déterminer sa qualité.

<p><b>Les Données</b></p> <p>Méfaits de l'école</p> <p>Les données ne sont pas données</p> <ul style="list-style-type: none"> <li>- pas évidentes,</li> <li>- pas disponibles aisément</li> <li>- pas gratuites.</li> </ul>	<p>A l'école on nous a appris à recevoir des « données ». Elles étaient indiscutables, intangibles, il fallait les utiliser toutes, ... si on était bon, elles permettaient d'atteindre un résultat. La vie nous a appris autre chose</p> <p>Les données sont des constructions, elles sont rarement évidentes, elles coûtent parfois très cher à produire.</p> <p>Elles donnent une vision du monde selon un certain point de vue. Elles ne sont pas le monde pas plus qu'une carte géographique n'est le territoire.</p>
---	--

<p><b>Cercle de la connaissance</b></p> <p><b>observation</b> : point de vue, acuité, ...</p> <p><b>Données</b> : nombreux choix techniques</p> <p><b>Calculs</b> : choix des méthodes et orga.</p> <p><b>Résultats</b> : choix des représentations</p> <p><b>Interprétation</b></p> <p>Connaissances</p>	<p>Le diagramme illustre le cycle de la connaissance. Au centre, un nuage est étiqueté 'Rien mêmes'. Des flèches forment un cycle : 'Observations' pointe vers 'Données', 'Données' pointe vers 'Calculs', 'Calculs' pointe vers 'Résultats', 'Résultats' pointe vers 'Interprétation', 'Interprétation' pointe vers 'Connaissances', et 'Connaissances' pointe vers 'Observations'.</p>
---	--

	À côté	Au-dessus	Dans	Au-dessous
Exemples	Médiamétrie, Google Analytics, affichage publicitaire	PageRank de Google, Digg, Wikipédia	Nombre d'amis Facebook, Retweet de Twitter, notes et avis	Recommandation Amazon, publicité comportementale
Données	Vues	Liens	Likes	Traces
Population	Échantillon représentatif	Vote censitaire, communautés	Réseau social, affinitaire, déclaratif	Comportements individuels implicites
Forme du calcul	Vote	Classements méritocratiques	Benchmark	Machine learning
Principe	Popularité	Autorité	Réputation	Prédiction

Cardon, p. 18

<p><b>Quelles données ?</b></p>	<p>Tous les algorithmes dont nous parlons (par exemple, ceux du tableau donné par la CNIL) ont besoin de données mais, selon le type, ces données seront de natures très différentes.</p> <p>Pour APB, le système d'information sur les élèves, les vœux des bacheliers, les préférences des formations.</p> <p>Pour un GPS, une cartographie détaillée et si possible à jour, pour prédire les crimes et délits, les bases de données de la police et de la justice, etc.</p> <p>Pour les systèmes prédictifs ou prescriptif internes à une entreprise ou à une administration, ce sont les bases de données internes. Ce qui va importer, c'est la taille de ces bases : par exemple, Amazone a suffisamment de clients pour faire des suggestions d'achats, Les Volcans ne pourront pas....</p>
<p>« Quelle étude faire » devient « quelles data » donner à la machine.</p>	<p>Ex : satisfaction et souhait des adhérents des RPC ? Entretiens, enquêtes (face à face, par courrier, comptages stat...) Observation de l'accès au site web - Accès à Facebook, Amazon, Google, médiathèque, ... EN....</p>
<p>Les données deviennent la ressource essentielle : d'où les slogans du type : l'or noir de demain</p> <p>Cf l'ouvrage : <b>L'empire des données</b> (Adrien Basdevant) à l'origine <b>d'un « coup data » !</b></p>	

<p><b>L'appariement de fichiers</b></p> <p>Ingénieurs sociaux, Société rationnelle</p>	<p>Le graal c'est de pouvoir rassembler sur une même personne des données venant de plusieurs sources différentes : santé (SNIIR), revenus, adresse, sexe et âge, ....C'était déjà l'idée, dès les années 30, d'un certain nombre « d'ingénieurs sociaux » dont l'exemple le plus notable est celui de René Carmille.</p> <p>René Carmille, après polytechnique, est contrôleur général de l'armement. Il s'intéresse aux fichiers de population utile à l'organisation de la conscription et de la mobilisation en cas de besoin. A ce titre il examine sur place l'utilisation des machines mécanographiques à cartes perforées qui sont en service dans les grandes entreprises américaines. Au retour, son projet est d'introduire cette technologie dans l'administration française : son grand projet c'est de rassembler derrière le même identifiant d'une personne toutes les informations dont disposent les différents services administratifs : démographiques, scolaire, justice, emploi, finance, ... Il propose son projet au gouvernement de front populaire qui n'en veut pas, plus tard il le propose au gouvernement de Vichy qui n'en veut pas davantage mais, comme il a été nommé directeur du Service National de la</p>
--	---

<p>Liberté individuelle, Informatique et liberté</p>	<p>Statistique (SNS), il impose dans le fichier de l'état civil un <b>identifiant à 13 chiffres</b> qui sera adopté par la sécurité sociale en 1946. A cette date, le numéro n'est utilisé que par l'INSEE qui a succédé au SNS (pour le répertoire national des personnes physiques et le fichier électoral) et la sécurité sociale pour sa gestion interne (maladie et retraite)<sup>3</sup>. On est donc loin du projet initial de Carmille de rassembler toutes les informations administratives derrière cet identifiant unique</p> <p>Ce répertoire national (nom, prénom, date de naissance, identifiant) étant sur cartes perforées, la question s'est posée fin des années 60 de passer de la mécanographie à l'informatique .... une technologie qui offrait de nouvelles possibilités. De façon non complètement élucidée, l'idée que ce fichier pourrait servir à rassembler des informations provenant de sources différentes fit son chemin<sup>4</sup>. Tant et si bien qu'en 1974, un article du journal « Le Monde » titrait : « SAFARI ou la chasse aux français » dans lequel son auteur, Philippe Bouchet, dénonçait le secret qui entourait le projet conjoint de l'INSEE et du ministère de l'intérieur de « fichier » l'ensemble de la population. L'affaire fait grand bruit, les débats sont nombreux qui aboutiront en 1978 à la Loi Informatique et libertés qui encadre très strictement les rapprochements de fichiers de personnes. En particulier, elle interdit de rapprocher des fichiers à des fins commerciales ou politiques.</p>
<p><b>Les cookies</b></p> <p><b>Cookies tiers</b></p>	<p><b>Il fallait donc que le marketing trouve autre chose....</b></p> <p>Quand je visite le site web philo63, quand je lis la page sur Montaigne, le site ne sait pas si j'ai d'abord été lire la page sur Erasme. Pour des sites qui font du e-commerce, ce serait très gênant. De ce fait, a été inventé un fichier (dénommé cookie) qui va, en particulier, contenir un « panier » dans lequel vont venir s'inscrire les achats effectués au fil des pages parcourues. Pour livrer la commande, le site de commerce va demander des informations de livraison et, quand le visiteur va quitter le site, le cookies va être récupéré par le commerçant. Et, comme le client va peut-être revenir, le cookies va être laissé sur l'ordinateur où il a été créé.</p> <p>En allant un peu plus loin, le commerçant ou media va essayer de « fidéliser » le client, lui faire ouvrir un « compte » en essayant de lui faire livrer des informations complémentaires.</p> <p>Où ça se complique c'est quand des sociétés, des régies publicitaires<sup>5</sup>, viennent installer leurs cookies (dits « cookies tiers ») pour en faire commerce. Par ce système, chaque site affilié sait ce qu'a fait chacun de ses clients sur tous les sites affiliés à la même régie. La régie publicitaire se constitue ainsi une banque de données rapidement géante (big data !). C'est par ce système qu'après avoir commandé une perceuse électrique sur un site de bricolage je vais retrouver une publicité pour perceuses électriques lorsque je vais consulter ma messagerie sur le portail d'orange. Et il y a des mécanismes encore plus « vicieux » du genre augmentation du prix lorsqu'un visiteurs revient sur le site qu'il avait</p>

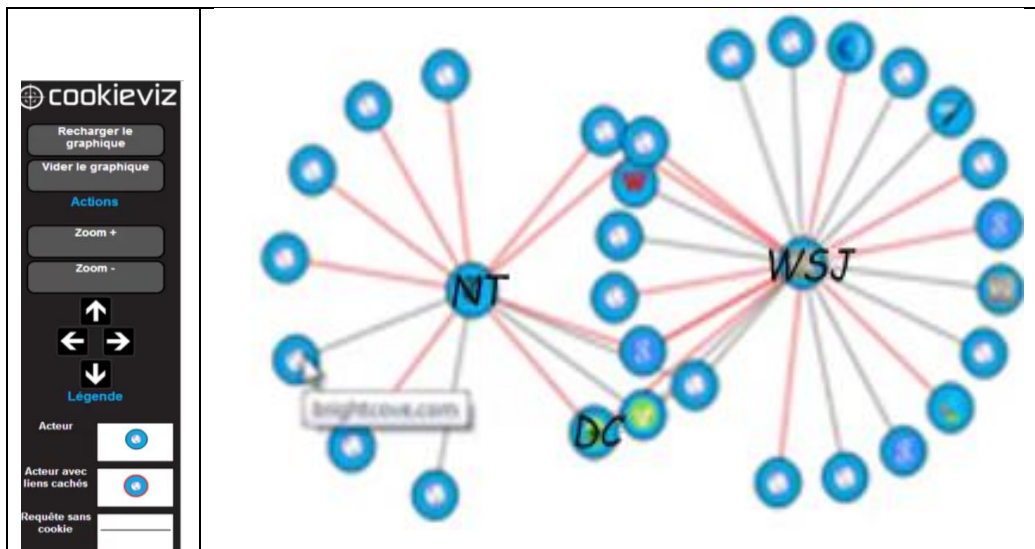
<sup>3</sup> Voir wikipedia pour en savoir plus sur René Carmille.

<sup>4</sup> Le projet a pris le nom de *Safari* ([Système automatisé pour les fichiers administratifs et le répertoire des individus](#)).

<sup>5</sup> Ces régies sont connues sous le terme de « ad-network » : weborama, double-click, critéo, ...



	<p>d'abord consulté avant de rechercher, par exemple, un meilleur tarif. <b>En fait, l'inventivité des gens de marketing paraît sans limites !</b></p> <p>Ces dispositifs sont le moyen de rassembler sur un même utilisateur des données qui proviennent de différentes sources : faute de pouvoir utiliser un identifiant commun style n°insee/SS (ce qui est interdit) l'appariement se fera sur un nom, un pseudo, un numéro de carte bleu, une adresse IP-, ...)<sup>6</sup>. Nos consommations en disent long sur nous-mêmes et, même si nous n'achetons rien sur Internet, la liste des sites et des pages que nous consultons en disent long sur nous-mêmes.</p>
--	--



Exemple : soit une visite du site du Wall Street Journal (WSJ sur la carte). Il est en relation avec une vingtaine d'autres sites pour leur signaler que l'on vient d'arriver chez lui.. Si, dans la même session on consulte le New-York Time (NT sur la carte), celui-ci, à son tour va signaler à une dizaine de sites que l'on vient d'arriver (dont quelques-uns communs au WSJ). Il y a un site commun particulier, celui de la régie publicitaire Double-Clic (noté DC sur la carte<sup>7</sup>) qui va communiquer à ses deux clients NT et WSJ (et à tous les autres intéressés) tous les cookies du visiteur.

Voir le site de la CNIL :

<https://linc.cnil.fr/fr/cookieviz-une-dataviz-en-temps-reel-du-tracking-de-votre-navigation>

<p><b>Les biais dans les données</b></p> <p><b>Biais =/ erreur</b></p>	<p>Une donnée peut être erronée : le bachelier a eu 18 en physique mais c'est 13 qui a été écrit dans la base de données. C'est une <b>erreur</b> de lecture, de saisie, de traitement, ... Elle est le plus souvent ponctuelle.</p> <p>Le <b>biais</b> désigne une erreur de méthode : penser que tout le monde se connecte sur Internet par exemple.</p>
--	--

<sup>6</sup> L'imagination des gens du marketing dans l'utilisation du web est impressionnante ... et il n'y a pas que du marketing commercial. Voir le dossier CNIL sur la publicité en ligne

[https://www.cnil.fr/sites/default/files/typo/document/Publicite\\_Ciblee\\_rapport\\_VD.pdf](https://www.cnil.fr/sites/default/files/typo/document/Publicite_Ciblee_rapport_VD.pdf)

<sup>7</sup> Racheté par Google en 2007.



<p><b>Les erreurs les plus graves imputées aux algorithmes proviennent de biais dans les données d'apprentissage</b></p>	<p>-Apprendre à traduire une langue à partir des textes de l'Union Européenne fait l'impasse sur l'existence de la première personne.</p> <p>- dans un système de recherche d'emploi, on s'aperçoit que le système propose des emplois moins bien payés aux femmes plutôt qu'aux hommes.</p> <p>-dans un système (US) de prescription de la liberté conditionnelle, les personnes de couleurs sont moins bien traitées que les blancs,</p> <p>-tag « gorille » attribué automatiquement à des photos d'identités de personnes noires, .....(rapport cnil p.32)</p> <p>Ce n'est pas le mécanisme d'apprentissage qui est défaillant mais bien les données que leurs concepteurs leur ont mis en apprentissage</p> <p>Cf.Cathy O'Neil : « Weapons of math destruction » (Armes de destruction mathématique)</p>
<p><b>La nécessité d'un prétraitement</b></p> <p>A supposer que les sources de data soient pertinentes, il faut des traitements préparatoires pour qu'elles soient utilisables.</p>	<p>Avant qu'un jeu de données soit utilisable par l'algorithme, il faut pratiquement toujours opérer des traitements de mise en forme, de tri, de codage, de correction, etc. Par exemple, des chercheurs travaillant sur le SNIIRAM<sup>8</sup> viennent de mettre en évidence une corrélation entre un médicament antidiabétique et l'apparition d'un cancer de la vessie. Après analyse par les médecins, il y a reconnaissance d'un lien de causalité.</p> <p>Ce type d'information qui n'avait pas été mise en évidence jusque-là peut avoir une grande importance et les résultats de l'étude incitent à exploiter ce SNIIRAM à grande échelle pour explorer des corrélations de ce type. Mais, comme cette base de données existe depuis près de 40 ans, on peut se demander pourquoi elle n'a pas été explorée plus tôt. La réponse donnée dans l'article relatant cette découverte (Le Monde du 30 janvier 2018) est que la base de données a été conçue pour gérer les remboursements de l'assurance maladie et pas pour faire de la recherche épidémiologique. <b>Pour cette recherche il a fallu un travail préparatoire de six ingénieurs pendant deux ans.</b></p> <p>Ces traitements préparatoires sont décidés par les responsables de chaque projet. Modifiant plus ou moins profondément le fichier d'apprentissage ils modifient les résultats fournis par l'algorithme. On peut vouloir « redresser » le fichier pour qu'il soit plus conforme à ce que l'on sait du phénomène décrit mais on peut aussi le « déformer » au profit de telle ou telle politique.</p>
<p><b>Les big-data</b> (mega données, données massives, ...)</p> <p>Clics, likes, pages, liens, ...</p>	<p>Au sens précis, il s'agit de méthodes (et donc de programmes) qui ont été conçues pour traiter les masses de données générées par l'utilisation d'Internet, des clics, des likes, des traces de navigation (pages et liens), ...des données souvent hétérogènes (textes, adresses, images, sons, vidéos, ...) le tout à traiter rapidement. Ceci a obligé à trouver de nouvelles façons de stocker</p>

<sup>8</sup> Le **Système national d'information inter-régimes de l'Assurance maladie (SNIIRAM)** est une base de données française maintenue par la [Caisse nationale de l'assurance maladie des travailleurs salariés](#). Il contient les données correspondant à plus d'un milliard de feuilles de soins par an.



