

Le « Deep learning », une révolution en Intelligence artificielle

Yann LECUN, chercheur et directeur du laboratoire de recherche en intelligence artificielle de Facebook, est nommé professeur invité sur la Chaire *Informatique et sciences numériques*

Une chaire créée en partenariat avec Inria

- Leçon Inaugurale le 04 février 2016 à 18h00 -

Fidèle à sa mission d'être toujours à la pointe des nouveaux développements et avancées scientifiques, le Collège de France a créé, il y a six ans en partenariat avec Inria, une chaire annuelle *Informatique et sciences numériques*¹. Chaque année s'y succèdent les plus grands chercheurs de domaines qui bouleversent nos quotidiens mais également le monde de la recherche et l'ensemble des sciences. Cette année, c'est Yann LeCun, l'un des plus éminents chercheurs en Intelligence artificielle, professeur à l'université de New York et directeur du laboratoire de recherche en intelligence artificielle de Facebook (FAIR), qui occupera cette chaire avec un cycle d'enseignement ouvert à tous.

Yann LeCun, spécialiste de l'apprentissage automatique des machines (« machine learning »), est l'un des pères du « Deep Learning » (« apprentissage profond »)²; une méthode à laquelle il se consacre depuis 30 ans, malgré le scepticisme qu'il rencontre au départ dans la communauté scientifique. Le Deep Learning, qui fait appel à la fois aux connaissances en neurosciences, aux mathématiques et aux progrès technologiques, est aujourd'hui plébiscité comme une véritable révolution dans le domaine de l'intelligence artificielle. Il a déjà permis d'immenses progrès et de multiples applications dans les domaines de la reconnaissance faciale et vocale, de l'étiquetage d'images, du traitement automatisé du langage ou encore de la vision par ordinateur.

« Les cerveaux humain et animal sont « profonds », dans le sens où chaque action est le résultat d'une longue chaîne de communications synaptiques (de nombreuses couches de traitement). Nous recherchons des algorithmes d'apprentissage correspondants à ces « architectures profondes ». Nous pensons que comprendre l'apprentissage profond ne nous servira pas uniquement à construire des machines plus intelligentes, mais nous aidera également à mieux comprendre l'intelligence humaine et ses mécanismes d'apprentissages », estime Yann LeCun.

Le Deep learning fait l'objet d'importants investissements privés, notamment de la part des grands acteurs du net, mais aussi publics. *« De plus en plus d'entreprises ont des masses de données gigantesques à exploiter, trier, indexer, et cela demande des ressources considérables. L'intelligence artificielle et le Deep learning peuvent aider à le faire de façon automatisée et plus efficace »,* confirme Yann LeCun qui reste prudent quant aux fantasmes que suscitent ces développements. *« De grands progrès ont été faits notamment en matière de reconnaissance visuelle et vocale - dans la reconnaissance*

automatique d'images, des réseaux neuronaux artificiels ont produit des algorithmes meilleurs que ceux conçus par des ingénieurs humains – mais nous sommes très loin de ce qu'un cerveau peut faire et nous n'en avons pas la prétention. Les animaux et les humains peuvent apprendre à voir, percevoir, agir et communiquer avec une efficacité qu'aucune machine ne peut approcher. D'autre part, il s'agit pour l'instant d'un apprentissage purement supervisé : on entraîne la machine à reconnaître l'image d'une voiture par exemple en lui montrant des milliers d'images et en la corrigeant quand elle fait erreur. Les humains découvrent le monde de façon non supervisée. L'apprentissage non supervisé est le défi scientifique auquel nous nous attelons. Tant que nous n'y serons pas parvenus, nous serons incapables de construire des systèmes intelligents”.

De la reconnaissance des tumeurs cancéreuses à la sécurité routière, les développements de l'intelligence artificielle et du Deep learning ouvrent de larges horizons que Yann LeCun exposera au Collège de France. Quant aux craintes générées par ces nouveaux domaines, il les comprend mais les modère, *“Même si un jour on construit des systèmes par certains aspects plus complexes ou performants que les humains, ils vont être construits pour des tâches spécifiques. On associe trop souvent l'intelligence artificielle aux qualités et aux défauts humains. Il n'y a aucune raison que les machines que l'homme concevra aient comme lui des désirs, des pulsions et des défauts ! »*

Yann LeCun donnera sa leçon inaugurale, *L'apprentissage profond : une révolution en intelligence artificielle*, le 4 février 2016. Ses cours auront lieu les vendredis à 14h30, à partir du 12 février. L'ensemble de son enseignement sera disponible sur le site de l'institution www.college-de-France.fr

¹ En 2012, Le Collège de France a aussi créé une chaire pérenne d'informatique : *Algorithmes, machines et langages*, avec pour titulaire le Pr Gérard Berry

² **Le Deep learning** est un ensemble de méthodes d'apprentissage automatique conçu sur la base de ce que l'on appelle parfois des réseaux de « neurones artificiels » à plusieurs couches (ou « convolutifs »). Ces réseaux sont capables de catégoriser les informations des plus simples aux plus complexes. Pour un objet par exemple, la première couche détecte des petits contours élémentaires, la seconde assemble ces contours en motifs puis les motifs en parties d'objets puis ces parties en objets.

Ces « neurones artificiels » n'ont rien de matériel. Ils sont en fait des fonctions mathématiques à plusieurs paramètres ajustables. Une phase d'apprentissage sur des objets connus permet de trouver les meilleurs paramètres en montrant par exemple à la machine des milliers d'images d'un chien, d'une voiture ou d'un sport ... L'un des enjeux étant de trouver des méthodes pour ajuster ces paramètres le plus rapidement et le plus efficacement possible.



Biographie

Yann LeCun est un chercheur en intelligence artificielle, apprentissage machine, vision artificielle et robotique. Il travaille depuis 30 ans sur «l'apprentissage automatique» (machine learning) et «l'apprentissage profond» (deep learning) : à savoir la capacité d'un ordinateur à reconnaître des représentations (images, textes, vidéos, sons) à force de les lui montrer de très nombreuses fois. Il a aussi contribué au développement de méthodes de compression d'image avec le format d'archivage DjVu.

Directeur depuis 2013 de *Facebook AI Research (FAIR)*, le centre de recherche en intelligence artificielle de Facebook, il est également professeur depuis 2003 à l'Université de New York (NYU), principalement affilié au *NYU Center for Data Science* et au *Courant Institute of Mathematical Sciences*.

Titulaire du diplôme de l'Ecole Supérieure d'Ingénieurs en Electronique et Electrotechnique (ESIEE) de Paris, il obtient un DEA de l'Université Pierre-et-Marie-Curie en 1984 et un doctorat en 1987. Après un post-doctorat à l'Université de Toronto, il rejoint les laboratoires AT&T Bell Laboratories à Holmdel (Etats-Unis) en 1988. Il prend la tête du département Traitement des Images d'AT&T Labs-Research en 1996 et devient professeur à NYU en 2003 après une brève période en tant que membre de la *NEC Research Institute* à Princeton. Il a fondé le NYU Center for Data Science en 2012 et en a été le directeur jusqu'à 2014.

Yann Le Cun est également le co-directeur du *Neural Computation and Adaptive Perception Program* du CIFAR et a co-dirigé pour NYU le *Moore-Sloan Data Science Environments Initiative*. Il a reçu le prix «IEEE Neural Network Pioneer Award» en 2014 et le «IEEE PAMI Distinguished Researcher Award» en 2015.

Il a publié près de 200 articles et documents sur l'apprentissage automatique, les réseaux neuronaux et la reconnaissance d'images, domaines dans lesquels il est considéré comme l'un des pionniers.



Présentation du cycle d'enseignement de Yann LeCun

Leçon inaugurale le 04 février 2016, à 18h00 :

L'apprentissage profond : une révolution en intelligence artificielle

Cours les vendredis à 14h30 (à partir du 12 février), suivis à 15h30 d'un séminaire en relation avec le cours :

- 12 février : Pourquoi l'apprentissage profond ?
- 19 février : Réseaux multi-couches et rétropropagation du gradient
- 26 février : L'apprentissage profond en pratique
- 04 mars : Réseaux convolutifs
- 25 mars : Réseaux convolutifs. Applications à la vision
- 01 avril : Réseaux récurrents. Applications au traitement du langage naturel.
- 08 avril : Raisonnement, attention, mémoire.

- 15 avril : L'Apprentissage non-supervisé (cours à 10h00)

Suivi d'un **colloque international toute la journée du 15 avril**

Programme complet : <http://www.college-de-france.fr/site/yann-lecun/index.htm>

(L'ensemble de l'enseignement de Yann LeCun sera disponible sur notre site à cette même adresse)

La chaire *Informatique et sciences numériques*

La chaire annuelle « *Informatique et sciences numériques* » a été créée en 2009 dans le cadre d'un partenariat entre le Collège de France et Inria. Elle a été inaugurée par Gérard Berry (nommé depuis professeur titulaire d'une chaire pérenne d'informatique, *Algorithmes, machines et langages*) et accueille chaque année un nouveau titulaire spécialiste reconnu d'un domaine (langages de programmation, sécurité informatique, Big Data, etc).

L'informatique, une science au cœur de notre quotidien et des enjeux scientifiques de demain

« *Le Collège de France a toujours eu pour vocation de soutenir les sciences nouvelles et de rester en alerte sur les évolutions scientifiques qui bouleversent nos sociétés ; la révolution informatique et numérique étant sans conteste l'une d'entre elles. Tous les pans de notre vie sont touchés par ces bouleversements, notre quotidien, notre économie mais aussi la recherche scientifique. L'informatique ouvre en effet pour de nombreuses autres sciences des horizons et des territoires de recherche jusque-là insoupçonnables, à travers notamment les nouvelles possibilités de calcul, de simulation et de modélisation. Le Collège de France ne pouvait l'ignorer* », a expliqué lors de sa création l'Assemblée des professeurs du Collège de France.

Pour les responsables Inria, « *La chaire annuelle Informatique et sciences numériques a marqué l'entrée de l'informatique en tant que discipline scientifique autonome au sein du Collège de France. C'est une étape importante dans la reconnaissance de notre domaine scientifique. Cette reconnaissance est fondamentale pour que l'informatique bénéficie d'un enseignement plus large. De ce dernier dépendra la capacité de notre pays à profiter pleinement des avancées offertes par les nouvelles technologies, tant sur le plan sociétal qu'économique.* »

Les précédents titulaires de cette chaire ont été :

Marie-Paule Cani (2014/2015) : Façonner l'imaginaire, de la création numérique 3D aux mondes virtuels animés

Nicholas Ayache (2013/2014) : Des images médicales au patient numérique

Bernard Chazelle (2012/2013) : L'algorithmique et les sciences

Serge Abiteboul (2011/2012) : Sciences des données : de la Logique du premier ordre à la Toile

Martin Abadi (2010/2011) : La Sécurité informatique.

Gérard Berry (2009/2010) : Penser, modéliser et maîtriser le calcul informatique.

À propos d'Inria

Inria, institut national de recherche dédié au numérique, promeut « l'excellence scientifique au service du transfert technologique et de la société ». Inria emploie 2700 collaborateurs issus des meilleures universités mondiales, qui relèvent les défis des sciences informatiques et mathématiques. Son modèle ouvert et agile lui permet d'explorer des voies originales avec ses partenaires industriels et académiques.

Inria répond ainsi efficacement aux enjeux pluridisciplinaires et applicatifs de la transition numérique et est à l'origine de nombreuses innovations créatrices de valeur et d'emplois. Inria est également très impliqué dans des actions en faveur de l'enseignement des sciences du numérique.

Inria et la recherche en intelligence artificielle :

La recherche sur l'intelligence artificielle a beaucoup apporté à l'informatique : des concepts fondamentaux comme les langages fonctionnels, les langages objets, le traitement de l'image, l'interface homme-machine, etc. C'est un axe très important pour Inria avec plusieurs équipes dans ce domaine, qui développent aussi bien des approches fondamentales de l'apprentissage que de nouvelles méthodes de reconnaissance d'images ou d'aide à la prise de décision. Les scientifiques Inria confrontent en permanence leurs recherches aux questions les plus actuelles et aux problèmes concrets que se posent aujourd'hui les entreprises du monde numérique et de nouveaux champs d'application comme la biologie ou la médecine. Certaines équipes collaborent avec le centre de recherche de Facebook (FAIR), récemment implanté en France et dirigé par Yann LeCun, sur des projets dans les domaines du traitement des images ou du langage, et des infrastructures logiques et physiques nécessaires aux systèmes d'intelligence artificielle. Ce type de collaboration est crucial pour tester des idées ou des algorithmes sur la complexité de données réelles et à grande échelle.

Contact presse Inria : Laurence Goussu / laurence.goussu@inria.fr - tel : 01 39 63 57 29

Les Enjeux de la Recherche en Intelligence Artificielle

Yann LeCun

Directeur, Facebook AI Research

Professeur, New York University

Chaire Informatique et Sciences Numériques
2015-2016

1 Qu'est-ce que l'Intelligence Artificielle

Qu'est-ce que l'intelligence ? Est-ce la capacité à percevoir le monde, à prédire le futur immédiat ou lointain, ou à planifier une série d'actions pour atteindre un but ? Est-ce la capacité d'apprendre, ou celle d'appliquer son savoir à bon escient ? La définition est difficile à cerner.

On pourrait dire que l'Intelligence Artificielle (IA) est un ensemble de techniques permettant à des machines d'accomplir des tâches et de résoudre des problèmes normalement réservés aux humains et à certains animaux.

Les tâches relevant de l'IA sont parfois très simples pour les humains, comme par exemple reconnaître et localiser les objets dans une image, planifier les mouvements d'un robot pour attraper un objet, ou conduire une voiture. Elles requièrent parfois de la planification complexe, comme par exemple pour jouer aux échecs ou au Go. Les tâches les plus compliquées requièrent beaucoup de connaissances et de sens commun, par exemple pour traduire un texte ou conduire un dialogue.

Depuis quelques années, on associe presque toujours l'intelligence aux capacités d'apprentissage. C'est grâce à l'apprentissage qu'un système intelligent capable d'exécuter une tâche peut améliorer ses performances avec l'expérience. C'est grâce à l'apprentissage qu'il pourra apprendre à exécuter de nouvelles tâches et acquérir de nouvelles compétences.

Le domaine de l'IA n'a pas toujours considéré l'apprentissage comme essentiel à l'intelligence. Dans le passé, construire un système intelligent consistait à écrire un programme "à la main" pour jouer aux échecs (par recherche arborescente), reconnaître des caractères imprimés (par comparaison avec des images prototypes), ou faire un diagnostic médical à partir des symptômes (par déduction logique à partir de règles écrites par des experts). Mais cette approche "manuelle" a ses limites.

2 L'apprentissage machine

Les méthodes manuelles se sont avérées très difficiles à appliquer pour des tâches en apparence très simples comme la reconnaissance d'objets dans les images ou la reconnaissance vocale. Les données venant du monde réel, les échantillons d'un son ou les pixels d'une image, sont complexes, variables et entachées de bruit.

Pour une machine, une image est un tableau de nombres indiquant la luminosité (ou la couleur) de chaque pixel, et un signal sonore une suite de nombres indiquant la pression de l'air à chaque instant.

Comment une machine peut-elle transcrire la suite de nombres d'un signal sonore en série de mots tout en ignorant le bruit ambiant, l'accent du locuteur et les particularités de sa voix ? Comment une machine peut-elle identifier un chien ou une chaise dans le tableau de nombres d'une image quand l'apparence d'un chien ou d'une chaise et des objets qui les entourent peuvent varier infiniment ?

Il est virtuellement impossible d'écrire un programme qui fonctionnera de manière robuste dans toutes les situations. C'est là qu'intervient l'apprentissage machine (qu'on appelle aussi apprentissage automatique). C'est l'apprentissage qui anime les systèmes de toutes les grandes compagnies de l'Internet. Elle l'utilisent depuis longtemps pour filtrer les contenus indésirables, ordonner des réponses à une recherche, faire des recommandations, ou sélectionner les informations intéressantes pour chaque utilisateur.

Un système entraînable peut être vu comme une boîte noire avec une entrée, par exemple une image, un son, ou un texte, et une sortie qui peut représenter la catégorie de l'objet dans l'image, le mot prononcé, ou le sujet dont parle le texte. On parle alors de systèmes de classification ou de reconnaissance des formes.

Dans sa forme la plus utilisée, l'apprentissage machine est *supervisé* : on montre en entrée de la machine une photo d'un objet, par exemple une voiture, et on lui donne la sortie désirée pour une voiture. Puis on lui montre la photo d'un chien avec la sortie désirée pour un chien. Après chaque exemple, la machine ajuste ses paramètres internes de manière à rapprocher sa sortie de la sortie désirée. Après avoir montré à la

machine des milliers ou des millions d'exemples étiquetés avec leur catégorie, la machine devient capable de classer correctement la plupart d'entre eux. Mais ce qui est plus intéressant, c'est qu'elle peut aussi classer correctement des images de voiture ou de chien qu'elle n'a jamais vues durant la phase d'apprentissage. C'est ce qu'on appelle la *capacité de généralisation*.

Jusqu'à récemment, les systèmes de reconnaissance des formes classiques étaient composés de deux blocs : un extracteur de caractéristiques (*feature extractor* en anglais), suivi d'un classifieur entraînable simple. L'extracteur de caractéristiques est programmé "à la main", et transforme le tableau de nombres représentant l'image en une série de nombres, un *vecteur de caractéristiques*, dont chacun indique la présence ou l'absence d'un motif simple dans l'image. Ce vecteur est envoyé au classifieur, dont un type commun est le *classifieur linéaire*. Ce dernier calcule une somme pondérée des caractéristiques : chaque nombre est multiplié par un poids (positif ou négatif) avant d'être sommé. Si la somme est supérieure à un seuil, la classe est reconnue. Les poids forment une sorte de "prototype" pour la classe à laquelle le vecteur de caractéristiques est comparé. Les poids sont différents pour les classifieurs de chaque catégorie, et ce sont eux qui sont modifiés lors de l'apprentissage. Les premières méthodes de classification linéaire entraînable datent de la fin des années 50 et sont toujours largement utilisées aujourd'hui. Elle prennent les doux noms de *perceptron* ou *régression logistique*.

3 Apprentissage profond et réseaux neuronaux

Le problème de l'approche classique de la reconnaissance des formes est qu'un bon extracteur de caractéristiques est très difficile à construire, et qu'il doit être repensé pour chaque nouvelle application.

C'est là qu'intervient l'apprentissage profond ou *deep learning* en anglais. C'est une classe de méthodes dont les principes sont connus depuis la fin des années 80, mais dont l'utilisation ne s'est vraiment généralisée que depuis environ 2012.

L'idée est très simple : le système entraînable est constitué d'une série de modules, chacun représentant une étape de traitement. Chaque module est entraînable, comportant des paramètres ajustables similaires aux poids des classifieurs linéaires. Le système est entraîné de bout en bout : à chaque exemple, tous les paramètres de tous les modules sont ajustés de manière à rapprocher la sortie produite par le système de la sortie désirée. Le qualificatif *profond* vient de l'arrangement de ces modules en couches successives.

Pour pouvoir entraîner le système de cette manière, il faut savoir dans quelle direction et de combien ajuster chaque paramètre de chaque module. Pour cela il faut calculer un *gradient*, c'est-à-dire pour chaque paramètre ajustable, la quantité par laquelle l'erreur en sortie augmentera ou diminuera lorsqu'on modifiera le paramètre d'une quantité donnée. Le calcul de ce gradient se fait par la méthode de *rétro-propagation*, pratiquée depuis le milieu des années 80.

Dans sa réalisation la plus commune, une architecture profonde peut être vue comme un réseau multi-couches d'éléments simples, similaires aux classifieurs linéaires, inter-connectés par des poids entraînaibles. C'est ce qu'on appelle un réseau neuronal multi-couches.

Pourquoi neuronal ? Un modèle extrêmement simplifié des neurones du cerveau les voit comme calculant une somme pondérée et activant leur sortie lorsque celle-ci dépasse un seuil. L'apprentissage modifie les efficacités des *synapses*, les poids des connexions entre neurones. Un réseau neuronal n'est pas un modèle précis des circuits du cerveau, mais est plutôt vu comme un modèle conceptuel ou fonctionnel. Le réseau neuronal est inspiré du cerveau un peu comme l'avion est inspiré de l'oiseau.

Ce qui fait l'avantage des architectures profondes, c'est leur *capacité d'apprendre à représenter le monde de manière hiérarchique*. Comme toutes les couches sont entraînaibles, nul besoin de construire un extracteur de caractéristiques à la main. L'entraînement s'en chargera. De plus, les premières couches

extraieront des caractéristiques simples (présence de contours) que les couches suivantes combineront pour former des concepts de plus en plus complexes et abstraits : assemblages de contours en motifs, de motifs en parties d'objets, de parties d'objets en objets, etc.

4 Réseaux Convolutifs, Réseaux Récurents

Une architecture profonde particulièrement répandue est le *réseau convolutif*. C'est un peu mon invention. J'ai développé les premières versions en 1988-1989 d'abord à l'Université de Toronto où j'étais postdoc avec Geoffrey Hinton (qui est maintenant à Google), puis aux Bell Laboratories, qui était à l'époque le prestigieux labo de recherche de la compagnie de télécommunication AT&T.

Les réseaux convolutifs sont une forme particulière de réseau neuronal multi-couches dont l'architecture des connexions est inspirée de celle du cortex visuel des mammifères. Par exemple, chaque élément n'est connecté qu'à un petit nombre d'éléments voisins dans la couche précédente. J'ai d'abord utilisé les réseaux convolutifs pour la reconnaissance de caractères. Mes collègues et moi avons développé un système automatique de lecture de chèques qui a été déployé largement dans le monde dès 1996, y compris en France dans les distributeurs de billets du Crédit Mutuel de Bretagne. A la fin des années 90, ce système lisait entre 10 et 20% de tous les chèques émis aux États-Unis. Mais ces méthodes étaient plutôt difficiles à mettre en oeuvre avec les ordinateurs de l'époque, et malgré ce succès, les réseaux convolutifs, et les réseaux neuronaux plus généralement, ont été délaissés par la communauté de la recherche entre 1997 et 2012.

En 2003, Geoffrey Hinton (de l'Université de Toronto), Yoshua Bengio (de l'Université de Montréal) et moi-même à NYU (l'Université de New York), décidions de démarrer un programme de recherche pour remettre au goût du jour les réseaux neuronaux, et pour améliorer leurs performances afin de raviver l'intérêt de la communauté. Ce programme, financé en partie par la fondation CIFAR (l'Institut Canadien de Recherches Avancées), je l'appelle parfois *la conspiration de l'apprentissage profond*.

En 2011-2012 trois événements ont soudainement changé la donne. Tout d'abord, les GPUs (*Graphical Processing Units*) capables de plus de mille milliards d'opérations par seconde sont devenus disponibles pour moins de 1000 Euros la carte. Ces puissants processeurs spécialisés, initialement conçus pour le rendu graphique des jeux vidéos, se sont avérés être très performants pour les calculs des réseaux neuronaux. Deuxièmement, des expériences menées simultanément à Microsoft, Google et IBM avec l'aide du laboratoire de Geoff Hinton, on montré que les réseaux profonds pouvaient diminuer de moitié les taux d'erreur des systèmes de reconnaissance vocale. Troisièmement, plusieurs records en reconnaissance d'image ont été battus par des réseaux convolutifs. L'évènement le plus marquant a été la victoire éclatante de l'équipe de Toronto dans la compétition de reconnaissance d'objets "ImageNet". Le diminution des taux d'erreurs était telle qu'une véritable révolution d'une rapidité inouïe s'est déroulée. Du jour au lendemain, la majorité des équipes de recherche en parole et en vision ont abandonné leurs méthodes préférées et sont passées aux réseaux convolutifs et autres réseaux neuronaux.

L'industrie de l'Internet a immédiatement saisi l'opportunité et a commencé à investir massivement dans des équipes de recherche et développement en apprentissage profond. L'apprentissage profond ouvre une porte vers des progrès significatifs en intelligence artificielle. C'est la cause première du récent renouveau d'intérêt pour l'IA.

Une autre classe d'architecture, les *réseaux récurrents* sont aussi revenus au goût du jour. Ces architectures sont particulièrement adaptées au traitement de signaux séquentiels, tels que le texte. Les progrès sont rapides, mais il y a encore du chemin à parcourir pour produire des systèmes de compréhension de texte et de traduction.

5 L'IA aujourd'hui. Ses enjeux

Les opportunités sont telles que l'IA, particulièrement l'apprentissage profond, sont vus comme des technologies d'importance stratégique pour l'avenir.

Les progrès en vision par ordinateur ouvrent la voie aux voitures sans chauffeur, et à des systèmes automatisés d'analyse d'imagerie médicale. D'ores et déjà, certaines voitures haut de gamme utilisent le système de vision de la compagnie israélienne MobilEye qui utilise un réseau convolutif pour l'assistance à la conduite. Des systèmes d'analyse d'images médicales détectent des mélanomes et autres tumeurs de manière plus fiable que des radiologues expérimentés. A Facebook, Google et Microsoft, des systèmes de reconnaissance d'images permettent la recherche et l'organisation des photos et le filtrage d'images violentes ou pornographiques.

Depuis plusieurs années déjà, tous les moteurs de reconnaissance vocale sur smartphone utilisent l'apprentissage profond.

Des efforts considérables de R&D sont consacrés au traitement du langage naturel : la compréhension de texte, les systèmes de question-réponse, les systèmes de dialogue pour les agents virtuels, et la traduction automatique. Dans ce domaine, la révolution de l'apprentissage profond a été annoncée, mais n'est pas encore achevée. Néanmoins, on assiste à des progrès rapides. Dans la dernière compétition internationale de traduction automatique, le gagnant utilisait un réseau récurrent.

6 La recherche en IA et les obstacles au progrès

Malgré tous ces progrès, nous sommes encore bien loin de produire des machines aussi intelligentes que l'humain, ni même aussi intelligentes qu'un rat.

Bien sûr, nous avons des systèmes qui peuvent conduire une voiture, jouer aux échecs et au Go, et accomplir d'autres tâches difficiles de manière plus fiable et rapide que la plupart des humains (sans parler des rats). Mais ces systèmes sont très spécialisés. Un gadget à 30 Euros nous bat à plate couture aux échecs, mais il ne peut rien faire d'autre.

Ce qui manque aussi aux machines, c'est la capacité à apprendre des tâches qui impliquent non seulement d'apprendre à représenter le monde, mais aussi à se remémorer, à raisonner, à prédire, et à planifier. Beaucoup de travaux actuels à Facebook AI Research et à DeepMind sont focalisés sur cette question. Une nouvelle classe de réseaux neuronaux, les *Memory-Augmented Recurrent Neural Nets* (réseaux récurrents à mémoire) est utilisée de manière expérimentale pour la traduction, la production de légendes pour les images, et les systèmes de dialogues.

Mais ce qui manque principalement aux machines, c'est le sens commun, et la capacité à *l'intelligence générale* qui permet d'acquérir de nouvelles compétences, quel qu'en soit le domaine. Mon opinion, qui n'est partagée que par certains de mes collègues, est que l'acquisition du sens commun passe par *l'apprentissage non-supervisé*.

Qu'il soit naturel ou artificiel, il y a trois formes principales d'apprentissage. Nous avons déjà parlé de l'apprentissage supervisé. Les deux autres formes sont l'apprentissage par renforcement, et l'apprentissage non supervisé.

L'apprentissage par renforcement désigne la situation où la machine ne reçoit qu'un simple signal, une sorte de récompense, indiquant si la réponse produite était correcte ou pas. Le scénario est similaire à l'entraînement d'un animal de cirque à qui l'on donne une friandise lorsqu'il exécute l'action désirée. Cette forme d'apprentissage nécessite de très nombreux essais, est utilisée principalement pour entraîner les machines à jouer à des jeux (par exemple les jeux vidéo ou le jeu de Go), ou à opérer dans des environnements simulés. On a assisté à un succès éclatant de l'apprentissage par renforcement combiné à l'apprentissage profond en la victoire récente

du programme de Go AlphaGo de DeepMind face au champion européen.

L'apprentissage non supervisé, quant à lui, est le mode principal d'apprentissage des animaux et des humains. C'est l'apprentissage que nous faisons par nous mêmes en observant le monde et en agissant. C'est en observant le monde que nous apprenons qu'il a trois dimensions, que des objets peuvent en cacher d'autres, que certains objets peuvent être déplacés, qu'un objet sans support tombe, qu'un objet ne peut pas être à deux endroits en même temps, etc.

C'est grâce à l'apprentissage non-supervisé que nous pouvons interpréter une phrase simple comme "Jean prend son portable et sort de la pièce". On peut inférer que Jean et son portable ne sont plus dans la pièce, que le portable en question est un téléphone, que Jean s'est levé, qu'il a étendu sa main pour attraper son portable, qu'il a marché vers la porte. Il n'a pas volé, il n'est pas passé à travers le mur. Nous pouvons faire cette inférence car nous savons comment le monde fonctionne. C'est le sens commun.

Comment acquérir ce sens commun ? Une hypothèse possible est *l'apprentissage prédictif*. Si l'on entraîne une machine à prédire le futur, elle ne peut y arriver qu'en élaborant une bonne représentation du monde et de ses contraintes physiques. Dans un scénario d'apprentissage prédictif, on montre à la machine un segment de vidéo, et on lui demande de prédire quelques images suivantes. Malheureusement, le futur est impossible à prédire exactement et la machine s'en tient à produire une image floue, une mixture de tous les futurs possibles.

Si l'intelligence est un gâteau au chocolat, le gâteau lui-même est l'apprentissage non-supervisé, le glaçage est l'apprentissage supervisé, et la cerise sur le gâteau est l'apprentissage par renforcement. Les chercheurs en IA sont dans la même situation embarrassante que les physiciens : 95% de la masse de l'univers est de nature complètement inconnue : matière noire et énergie noire. La matière noire de l'AI est la génoise au chocolat de l'apprentissage non-supervisé.

Tant que le problème de l'apprentissage non-supervisé ne sera pas résolu, nous n'aurons pas de machines vraiment intelligentes. C'est une question fondamentale scientifique et mathématique, pas une question de technologie. Résoudre ce problème pourra prendre de nombreuses années ou plusieurs décennies. A la vérité, nous n'en savons rien.

7 A quoi ressembleront les machines intelligentes de demain ?

Si nous arrivons à concevoir des techniques d'apprentissage machine aussi générales et performantes que celles de la nature, à quoi ressembleront les machines intelligentes de demain ?

Il est très difficile de nous imaginer une entité intelligente qui n'ait pas toutes les qualités et les défauts des humains car l'humain est notre seul exemple d'entité intelligente. Comme tous les animaux, les humains ont des pulsions et des instincts gravés dans notre cerveau reptilien par l'évolution pour la survie de l'espèce. Nous avons l'instinct de préservation, nous pouvons devenir violents lorsque nous sommes menacés, nous désirons l'accès aux ressources pour ne pas mourir de faim, ce qui peut nous rendre jaloux. Nos instincts d'animaux sociaux nous conduisent aussi à rechercher la compagnie d'autres humains. Mais les machines intelligentes n'auront aucune raison de posséder ces pulsions et instincts. Pour qu'elles les aient, il faudrait que leurs concepteurs les construisent explicitement.

Les machines intelligentes du futur auront des sentiments, des plaisirs, des peurs, et des valeurs morales. Ces valeurs seront une combinaison de comportements, d'instincts et de pulsions programmées avec des comportements appris.

Dans quelques décennies, quand nous pourrons peut-être penser à concevoir des machines réellement intelligentes, nous devons répondre à la question de comment aligner les valeurs des machines avec les valeurs morales humaines.

Mais c'est un futur lointain où l'on pourra donner de l'autonomie aux machines. D'ici là, les machines seront certes intelligentes, mais pas autonomes. Elles ne seront pas à même de définir leurs propres buts et motivations. L'ordinateur de votre voiture s'en tiendra à conduire votre voiture en toute sécurité. L'IA sera un amplificateur de notre intelligence, et non un substitut pour celle-ci.

Malgré les déclarations de certaines personnalités, un scénario à la Terminator est immensément improbable. Tout d'abord, il faut garder à l'esprit que l'apparition de l'IA ne sera pas un événement singulier, ni le fait d'un groupe isolé. Le progrès de l'IA sera progressif et ouvert. Comprendre l'intelligence est une des grandes questions scientifiques de notre temps. Aucune organisation, si puissante soit-elle, ne peut résoudre ce problème en isolation. La conception de machines intelligentes nécessitera la collaboration ouverte de la communauté de la recherche entière.

8 Faut-il avoir peur de l'IA ?

L'IA n'éliminera donc pas l'humanité de sa propre initiative.

Mais comme toute technologie puissante, l'IA peut être utilisée pour le bénéfice de l'humanité entière ou pour le bénéfice d'un petit nombre aux dépens du plus grand nombre.

L'émergence de l'AI va sans doute déplacer des métiers. Mais elle va aussi sauver des vies (par la sécurité routière et la médecine). Elle va très probablement s'accompagner d'une croissance de la production de richesses par habitant. La question pour les instances dirigeantes est comment distribuer ces nouvelles richesses, et comment former les travailleurs déplacés aux nouveaux métiers créés par le progrès technologique. C'est une question politique et non technologique. C'est une question qui n'est pas nouvelle: l'effet du progrès technologique sur le marché du travail existe depuis la révolution industrielle. L'émergence de l'IA n'est qu'un symptôme de l'accélération du progrès technologique.