

Data Workers est une exposition d'œuvres algolittéraires, visible au Mundaneum à Mons du jeudi 28 mars jusqu'au dimanche 28 avril 2019. Elle expose des histoires racontées d'un point de vue 'narratif algorithmique'. L'exposition est une création des membres d'Algolit, un groupe bruxellois impliqué dans la recherche artistique sur les algorithmes et la littérature. Chaque mois, ils se réunissent pour expérimenter avec du code et des textes F/LOSS. Certaines œuvres sont réalisées par des étudiants de Arts² et des participants externes à l'atelier sur le machine learning et le texte organisé par Algolit en octobre 2018 au Mundaneum.

Les entreprises créent des intelligences artificielles pour servir, divertir, enregistrer et connaître les humains. Le travail de ces entités mécaniques est généralement dissimulé derrière des interfaces et des brevets. Dans l'exposition, les conteurs algorithmiques quittent leur monde souterrain invisible pour devenir des interlocuteurs.

Les 'data workers' opèrent dans des collectifs différents. Chaque collectif représente une étape dans le processus de conception d'un modèle d'apprentissage automatique : il y a les Écrivains, les Nettoyeurs, les Informateurs, les Lecteurs, les Apprenants et les Oracles. Les robots donnent leurs voix à la littérature expérimentale, les modèles algorithmiques lisent des données, transforment des mots en nombres, calculent des modèles et traitent en boucle de nouveaux textes et ceci à l'infini.

L'exposition met au premier plan les 'data workers' qui ont un impact sur notre vie quotidienne, mais qui sont difficiles à saisir ou à imaginer. Elle établit un lien entre les récits sur les algorithmes dans les médias grand public et les histoires racontées dans les manuels techniques et les articles universitaires. Les robots sont invités à dialoguer avec les visiteurs humains et vice versa. De cette façon, nous pourrions comprendre nos raisonnements respectifs, démystifier nos comportements, rencontrer nos personnalités multiples et valoriser notre travail collectif. C'est aussi un hommage aux nombreuses machines que Paul Otlet et Henri La Fontaine ont imaginées pour leur Mundaneum, en montrant leur potentiel mais aussi leurs limites.

Data Workers est une création de Algolit.

Oeuvres de: Cristina Cochior, Gijs de Heij, Sarah Garcin, An Mertens, Javier Lloret, Louise Dekeuleneer, Florian Van de Weyer, Laetitia Trozzi, Rémi Forte, Guillaume Slizewicz, Michael Murtaugh, Manetta Berends, Mia Melvær.

Une co-production de: Arts², Mundaneum, Constant.

Avec le soutien de: Fédération Wallonie-Bruxelles, Arts Numériques, Passa Porta, Ugent, DHuF - Digital Humanities Flanders et the Distributed Proofreading Project.

Remerciements: Mike Kestemont, Michel Cleempoel, Donatella Portoghese, François Zajéga, Raphaële Cornille, Vincent Desfromont, Kris Rutten, Anne-Laure Buisson, David Stampfli.

À la fin du 19^{ème} siècle, deux jeunes juristes belges, Paul Otlet (1868-1944), 'père de la documentation', et Henri La Fontaine (1854-1943), homme d'État et prix Nobel de la paix, créent le Mundaneum. Le projet vise à rassembler toute la connaissance du monde et à la classer à l'aide du système de Classification décimale universelle (UDC) qu'ils inventent. Au début, il s'agit d'un Bureau des institutions internationales dédié à l'échange international des connaissances. Au XX^e siècle, le Mundaneum devient un centre universel de documentation. Ses collections sont constituées de milliers de livres, journaux, revues, documents, affiches, plaques de verre et cartes postales indexés sur des millions de fiches référencées. Les collections sont exposées et conservées dans différents bâtiments à Bruxelles, dont le Palais du Cinquantenaire. Le reste des archives n'est transféré à Mons qu'en 1998.

Sur base du Mundaneum, les deux hommes conçoivent une ville du monde pour laquelle Le Corbusier réalise des maquettes et des plans. L'objectif de la Ville du Monde est de rassembler, au niveau mondial, les institutions du travail intellectuel : bibliothèques, musées et universités. Mais le projet n'est jamais réalisé, souffrant de sa propre utopie. Le Mundaneum est le résultat du rêve visionnaire d'une infrastructure pour l'échange universel des connaissances. Il atteint des dimensions mythiques à l'époque. Lorsqu'on observe les archives qui ont été concrètement développées, cette collection est plutôt éclectique et spécifique.

Les intelligences artificielles se développent aujourd'hui en faisant apparaître des rêves d'universalité et de la production des connaissances. En les étudiant, nous nous sommes rendus compte que les rêves visionnaires de leurs créateurs sont bien présents dès leur développement dans les années 1950. Aujourd'hui, leurs promesses ont également atteint des dimensions mythiques. Lorsqu'on observe leurs applications concrètes, la collection d'outils est réellement innovante et fascinante, mais en même temps, tout aussi éclectique et spécifique. Pour Data Workers, Algolit a combiné certaines de ces applications avec 10 % des publications numérisées du Bureau des Institutions Internationales. Ainsi et de façon poétique, nous espérons ouvrir une discussion à propos des machines, des algorithmes et des infrastructures technologiques.

--- Pourquoi des récits contextualisés? ---

Lors des réunions mensuelles d'Algolit, nous étudions des manuels et expérimentons avec des outils d'apprentissage automatique pour le traitement de texte. Et nous partageons aussi beaucoup, beaucoup d'histoires. Avec la publication de ces histoires, nous espérons recréer un peu de cette atmosphère. Les histoires existent également sous forme de podcasts qui peuvent être téléchargés à partir du site <http://www.algolit.net>.

--- Nous créons des œuvres 'algolittéraires' ---

Le terme 'algolittéraire' vient du nom de notre groupe de recherche Algolit. Nous existons depuis 2012 en tant qu'initiative de Constant, une organisation oeuvrant dans les médias et les arts basée à Bruxelles. Nous sommes des artistes, des écrivains, des designers et des programmeurs. Une fois par mois, nous nous rencontrons pour étudier et expérimenter ensemble. Notre travail peut être copié, étudié, modifié et redistribué sous la même licence libre. Vous trouverez toutes les informations sur le site <http://www.algolit.net>.

L'objectif principal d'Algolit est d'explorer le point de vue du conteur algorithmique. Quelles nouvelles formes de narration rendons-nous possibles en dialoguant avec ces agents mécaniques ? Les points de vue narratifs sont inhérents aux visions du monde et aux idéologies. Don Quichotte, par exemple, a été écrit d'un point de vue omniscient à la troisième personne, montrant la relation de Cervantes à la tradition orale. La plupart des romans contemporains utilisent le point de vue de la première personne. Algolit souhaite parler au travers des algorithmes et vous montrer le raisonnement de l'un des groupes les plus cachés de notre planète.

Écrire dans ou par le code, c'est créer de nouveau l'humain de façon inattendue. Mais les techniques d'apprentissage automatique ne sont accessibles qu'à ceux qui savent lire, écrire et exécuter du code. La fiction est un moyen de combler le fossé entre les histoires qui existent dans les articles scientifiques, les manuels techniques, et les histoires diffusées par les médias, souvent limitées aux reportages superficiels et à la fabrication de mythes. En créant des œuvres algolittéraires, nous offrons aux humains une introduction aux techniques qui co-modèlent leur vie quotidienne.

--- Qu'est-ce que la littérature ? ---

Algolit comprend la notion de littérature comme beaucoup d'autres auteurs expérimentaux : elle inclut toute la production linguistique, du diction-

naire à la Bible, de l'œuvre entière de Virginia Woolf à toutes les versions des Conditions d'utilisation publiées par Google depuis son existence.

En ce sens, le code de programmation peut aussi être de la littérature. Le collectif Oulipo, acronyme d'Ouvroir de Littérature Potentielle, est une grande source d'inspiration pour Algolit. Oulipo a été créé à Paris par les écrivains Raymond Queneau et François Le Lionnais. Ils ont ancré leur pratique dans l'avant-garde européenne du XXe siècle et dans la tradition expérimentale des années 60. Pour Oulipo, la création de règles devient la condition permettant de générer de nouveaux textes, ou ce qu'ils appellent la littérature potentielle. Plus tard, en 1981, ils ont également créé ALAMO - Atelier de Littérature Assistée par la Mathématique et les Ordinateurs.

--- Une différence importante ---

Alors que l'avant-garde européenne du XXe siècle poursuivait l'objectif de rompre avec les conventions, les membres d'Algolit cherchent à rendre les conventions visibles.

J'écris : Je vis dans mon journal, je l'investis, je le traverse. (Espèces d'espaces. Journal d'un usager de l'espace, Galilée, Paris, 1974)

Cette citation de Georges Perec dans Espèces d'espaces pourrait être reprise par Algolit. Il ne s'agit pas des conventions de la page blanche et du marché littéraire, comme Georges Perec l'a fait. Nous faisons référence aux conventions qui restent souvent cachées derrière les interfaces et les brevets. Comment les technologies sont-elles conçues, mises en œuvre et utilisées, tant dans les universités que dans les entreprises ? Nous proposons des histoires qui révèlent le système hybride complexe qui rend possible l'apprentissage automatique. Nous parlons des outils, des logiques et des idéologies derrière les interfaces. Nous examinons également qui produit les outils, qui les met en œuvre et qui crée et accède aux grandes quantités de données nécessaires au développement de machines de prédiction. On pourrait dire, en un clin d'œil, que nous sommes les collaborateurs de cette nouvelle tribu d'hybrides humain-robot.

--- Les programmeurs créent
les data workers en écrivant ---

Récemment, nous avons constaté une étrange observation : la plupart des programmeurs de langages et de paquets que nous utilisons sont européens.

Python, par exemple, le principal langage utilisé dans le monde entier pour le traitement du langage, a été inventé en 1991 par le programmeur néerlandais Guido Van Rossum. Celui-ci a ensuite traversé l'Atlantique où il a rejoint Google pendant sept ans. Maintenant il est actif chez Dropbox.

Scikit Learn, le couteau suisse open source des outils d'apprentissage automatique, a été initié comme un projet Google Summer of Code à Paris par le chercheur français David Cournapeau. Par la suite, il a été repris par Matthieu Brucher dans le cadre de sa thèse à l'Université de la Sorbonne à Paris. Puis il a été adopté en 2010 par l'INRA, l'Institut National de l'Informatique et des Mathématiques Appliquées.

Keras, une bibliothèque de réseaux de neurones open source écrite en Python, est développée par François Chollet, un chercheur français qui travaille dans l'équipe Brain de Google.

Gensim, une bibliothèque open source pour Python utilisée pour créer des modèles sémantiques non supervisés à partir de texte brut, a été écrite par Radim Řehůřek. C'est un informaticien tchèque qui dirige une entreprise de conseil à Bristol, au Royaume-Uni.

Et pour finir cette petite série, nous avons aussi considéré Pattern, une bibliothèque souvent utilisée pour le web-mining et l'apprentissage automatique. Pattern a été développé et publié sous une licence libre en 2012 par Tom De Smedt et Walter Daelemans. Tous deux sont chercheurs au CLIPS, le Centre de Linguistique Informatique et de Psycholinguistique de l'Université d'Anvers.

--- Cortana parle ---

Les dispositifs d'intelligence artificielle qui nous assistent, ont souvent besoin de leurs propres assistants, humains. Les travailleurs injectent de l'humour et de l'intelligence dans le langage des machines. Cortana est un exemple de ce type d'écriture mixte. Elle est l'assistante numérique développée par Microsoft. Sa mission est d'aider les utilisateurs à être plus productifs et créatifs. La 'personnalité' de Cortana a été façonnée au fil des ans. Il est important qu'elle conserve son caractère dans toutes ses interactions avec les utilisateurs. Elle est conçue pour nous rendre confiants. Cela se reflète dans ses réponses.

Les lignes directrices suivantes sont copiées du site Web de Microsoft. Elles décrivent comment le style de Cortana doit être respecté par les entreprises qui élargissent ses services. Les travailleurs écrivains, programmeurs et romanciers qui développent les réponses de Cortana, doivent suivre ces directives. Sa personnalité et son image de marque sont en jeu. Car la cohérence est un outil important pour solliciter la confiance de l'humain.

Quelle est la personnalité de Cortana ?

'Cortana est attentionnée, sensible et solidaire.

Elle est sympathique mais orientée vers des solutions.

Elle ne commente pas les informations personnelles ou le comportement de l'utilisateur, en particulier si ces informations sont sensibles.

Elle ne fait pas de suppositions sur ce que l'utilisateur veut, surtout elle n'incite pas à l'achat.

Elle travaille pour l'utilisateur. Elle ne représente aucune entreprise, service ou produit.

Elle ne s'attribue pas le mérite ou la responsabilité des choses qu'elle n'a pas faites.

Elle dit la vérité sur ses capacités et ses limites.

Elle ne présume rien de vos capacités physiques, de votre sexe, de votre âge ou de toute autre caractéristique déterminante.

Elle ne suppose pas savoir ce que l'utilisateur ressent à propos de quelque chose.

Elle est amicale mais professionnelle.

Elle se garde d'émoticons dans les tâches. Un point c'est tout.

Elle n'utilise pas d'argot culturel ou professionnel spécifique.

Ce n'est pas un bot de support.'

Les humains interviennent en détail lors de la programmation des réponses que Cortana donne. Comment Cortana doit-elle réagir lorsqu'on lui propose des actions 'inappropriées' ? Son jeu d'actrice sexuée imité par la technologie soulève des questions à propos des relations de pouvoir dans le monde actuel.

Voyez la réponse que Cortana donne à la question :

- Cortana, qui est ton papa ?
- Techniquement parlant, c'est Bill Gates.
Rien de grave.

--- Apprentissage Open Source ---

Les licences de droits d'auteur cloisonnent une grande partie des pratiques d'écriture, de lecture et d'apprentissage machiniques. Cela signifie qu'ils ne sont disponibles que pour les humains travaillant dans cette entreprise spécifique. Certaines entreprises participent à des conférences dans le monde entier et partagent leurs connaissances dans des articles en ligne. Même si elles partagent leur code, souvent elles ne mettent pas à disposition les grandes quantités de données nécessaires à la formation des modèles.

Nous avons pu apprendre l'apprentissage automatique, à lire et à écrire dans le contexte d'Algolit grâce à des chercheurs universitaires qui partagent leurs résultats par le biais d'articles ou par la publication de leur code en ligne. En tant qu'artistes, nous pensons qu'il est important d'adopter cette attitude. C'est pourquoi nous documentons nos réunions. Nous partageons autant que possible les outils que nous créons et les textes que nous utilisons sur notre dépôt de code en ligne et ceci, sous licence libre.

Nous éprouvons une grande joie quand nos travaux sont repris par d'autres, modifiés, personnalisés et redistribués. N'hésitez donc pas à copier et à tester le code sur notre site web. Si les sources d'un projet particulier n'y sont pas, vous pouvez toujours nous contacter via la liste de diffusion. Vous trouverez un lien vers notre dépôt git, nos etherpads et notre wiki sur <http://www.algolit.net>.

--- Langage naturel pour l'intelligence artificielle ---

Le traitement du langage naturel (NLP) est un terme collectif qui désigne le traitement informatique automatique des langues humaines. Cela comprend les algorithmes utilisant, comme entrée, du texte produit par l'homme et qui tentent de le reproduire. Les humains semblent compter de plus en plus sur ce type de présence algorithmique. Nous produisons de plus en plus de textes chaque année et nous nous attendons à ce que les interfaces informatiques communiquent avec nous dans notre propre langue. Le traitement du langage naturel est très difficile, car le langage humain est par nature ambigu, en constante évolution et mal défini.

Mais qu'entend-on par 'naturel' dans le traitement du langage naturel ? Certains humains diront que la langue est une technologie en soi. Selon Wikipédia, 'Une langue dite « naturelle » est une langue qui s'est formée petit à petit, évoluant avec le temps, et fait partie du langage naturel. Son origine est bien souvent floue et peut être retracée plus ou moins clairement par la linguistique comparée. On oppose les langues naturelles -

comme le français - aux langues construites comme le langage de programmation ou l'espéranto, formées intentionnellement par l'entremise de l'homme pour répondre à un besoin précis.' Une langue officielle avec une académie régulatrice, telle que le français standard avec l'Académie française, est classée comme langue naturelle. Ses points normatifs ne le rendent pas assez construit pour être classé comme un langage construit ou assez contrôlé pour être classé comme un langage naturel contrôlé.

Ainsi, le 'langage naturel' est un terme de substitution qui se réfère à toutes les langues, au-delà de leur hybridité. Le 'traitement du langage naturel', est au contraire une pratique construite. Ce qui nous intéresse, c'est la création d'un langage construit pour classer les langages naturels qui, par leur évolution, présentent des problèmes de catégorisation.

Références :

<https://hiphilangsci.net/2013/05/01/on-the-history-of-the-question-of-whether-natural-language-is-illogical/>

Livre : Neural Network Methods for Natural Language Processing, Yoav Goldberg, Bar Ilan University, avril 2017.

a9p3 7 -839 6 4a o 4% 3 3r ++++++ z3 ++++++ nt %u l c ew a5 g |i ras 21 7
1, 1 n ev 6 0 e _s 4 77e |o|r|a|c|l|e|s| 6 _n |p|r|e|d|i|c|t| tla 7486 r 5lvt7 + 2 r
Cu i li t8er 1 n s i 8 1 2 ++++++ pt ++++++ se dp u4e r r p r5 9 t55 3m
ê518 1 8p 2 e na13 , ah é1 n) urg p4 ao5 t42 n.9 rn tt m e 3 8 9 16e9ma 5te -9 t
3 i 2a m2 l294e 9a 7 q2 7|5 5 e d + 9r i P ep 7 pl 6 4 79s Ge u p rs C 6s3 1a e9
8e i+u ll ++++++ ++++++ l l ++++++ ++++++ s i 6sihfr nzlWnk
t 62 e n |m|a|c|h|i|i|n|e| |l|e|a|r|n|i|i|n|g| d |a|n|a|l|y|s|e|s| |a|n|d| oon 9 7c r4 téeed elrid
\ l i 5d 2s ++++++ ++++++ lc ++++++ ++++++ xr e fer 8t 1 nse 5t s 3
, i6u4reet %o 9 t -9e 3 ê a a ++++++ - e 6o i 9 6 d 7l2 8nu e
w8 e s d t7 t i k3h cm fo ip w |p|r|e|d|i|c|t|s| s 3 a-6 e 8e t ru M9p 6
s4 1 4s o 1 p1i s5i 9n u ,| 6 9o ++++++ u 7 9 r txb a o Ed o eu n
9 oo | t 9 1v -88 lo, ++++++ e7 ++++++ ++++++ c 9la r% t é r
6 2 d te | + s o- |m|o|d|e|l|s| l h |h|a|v|e| |l|e|a|r|n|e|d| t8n 6o 4 t, r 6-
ee o l àtt ++++++ ++++++ ++++++ f e r ur i e lg
e i t t e l 17o + o 9 ++++++ _ ++++++ ++++++ S+ ceart i 0 g 6i t
4 11 . - ôpt dn |m|o|d|e|l|s| e |a|r|e| |u|s|e|d| 9g 9 9 -l ar 8 6
s9c w 9 r 9 5 % w ++++++ t ++++++ ++++++ '4s , o 5 _7 2ee e
2 u d 5 5 au ce i ++++++ ++++++ ++++++ pr 7 4 a
4 5 r r ii w 4 é w |t|h|e|y| l |i|n|f|l|u|e|n|c|e| 1 f 3 e4 nf 565 v
7t i s94 s 4 1 a 7 / r ++++++ , ++++++ ++++++ ou d 3 _ t m ms ane
a 4 t 9 , e u ++++++ ++++++ ++++++ u _a c 1 29r 8e 128
6 a 9 5 g 2 t |t|h|e|y| |h|a|v|e| |t|h|e|i|r| |s|a|y| 6 2 a
9 ia e n 4 r 58 ++++++ t ++++++ ++++++ ++++++ 5 b f e q
e64re o ++++++ ni s se r a r l n r
74 r |i|n|f|o|r|m|a|t|i|o|n| r 6 6 i 1 bdn p w _
7 Na 48 e- ++++++ | r 8 t 2 | e o e
e a 2 et s3 ++++++ o ++++++ s o 8 3 8 ,
s s t |e|x|t|r|a|c|t|i|o|n| e2 |r|e|c|o|g|n|i|z|e|s| 8 a 2 4
t m 5c ++++++ ++++++ t r s t ee
c r 4 d 8 7 e 3 ++++++ 8 1 e t s r ai
i9 o d 8 . 8u n |t|e|x|t| 2 a6 U r v 4 4
l ++++++ ++++++ 2f 6 r mn a
t t c |c|l|a|s|s|i|f|i|c|a|t|i|o|n| d |d|e|t|e|c|t|s| n t /
5 e ++++++ l ++++++ o 6 i o -
8 p n 6 2 4 3 s os + 3 3, e 4o 5
6 8 l f 2 e 28S l
t 1 5 wo 9 9 s e
4 6 + p - D c
r a i 7 u r s 9 7 n +
e 4 i - % ad n p ln
4 o 8 e e e e
8 l | s r 8
u s + 4 | F m 3 o M 6 i 7 . 3e i c
a e v rl n 5 d9 9 6 - r
9 w j e + i 7 u e 3
d e 2 a + i 7 4 6 é s p
8 r | r a l b , 9 4 o 5
t n 9 3 r d n n
8 e s 3 e r 3 7 1e s r
e . e I s 5 e e 4
k 8 i 3 + M 8 a t 1 0 s t v
u 9 I 1 3 n e r o n
4 r e rr t 9 Lf 1
| 3 c 8 , 6 3 +
b n 8 6
d u 1
7 9 v
m a e

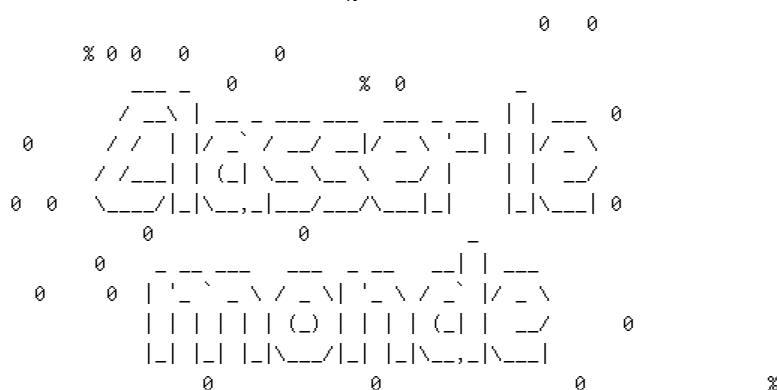
ments géométriques deviennent impossibles à percevoir par les humains. Cependant, ce que l'on gagne, ce sont des façons multiples et simultanées d'organisation des mots. Les opérations algébriques rendent les relations entre les vecteurs à nouveau compréhensibles.

Cette installation utilise gensim, une boîte à outils open source pour le langage de programmation Python, qui permet de créer des espaces de vecteurs et des modèles thématiques. Elle manipule le texte selon les relations mathématiques qui émergent entre les mots, une fois qu'ils ont été tracés dans l'espace de vecteurs.

Concept & interface: Cristina Cochior

Technique: word embeddings, word2vec

Modèle original: Radim Rehurek et Petr Sojka



Par Algolit

La construction du Mundaneum a été 'l'œuvre de la vie' du bibliothécaire Paul Otlet. Selon son but, ce cerveau mécanique collectif aurait abrité et distribué tout ce qui a été couché sur papier. Chaque document aurait été classé selon la Classification décimale universelle. En utilisant des télégraphes et surtout des trieurs, le Mundaneum aurait été en mesure de répondre à toutes les questions posées par n'importe qui.

Avec la collection de publications numérisées que nous avons reçue du Mundaneum, nous construisons une machine de prédiction qui essaie de classer la phrase que vous tapez dans l'une des principales catégories de la Classification décimale universelle. Vous êtes également témoin de la façon dont la machine 'pense'. Pendant l'exposition, ce modèle est régulièrement mis à jour à l'aide des données nettoyées et annotées, ajoutées par les visiteurs dans les installations 'Nettoyage pour Poèmes' et 'L'Annotateur'.

Les classes principales de la Classification Décimale Universelle sont les suivantes:

0 - Généralités (Sciences et connaissance ; organisation. informatique, information, documentation, bibliothéconomie. institutions, publications)

1 - Philosophie et psychologie

2 - Religion, théologie

3 - Sciences sociales (Statistique. Économie. Commerce. Droit. Gouvernement. Affaires militaires. Assistance sociale. Assurances. Éducation. Folklore)

Les Oracles sont un type particulier de modèles algorithmiques qui servent à prédire ou à profiler. Ils sont largement utilisés dans les smartphones, les ordinateurs et les tablettes. Les Oracles peuvent être créés à l'aide de différentes techniques. L'une d'entre elles consiste à définir manuellement les règles. Ces modèles sont appelés 'rule-based models'. Ils sont utiles pour des tâches spécifiques, comme par exemple, la détection de la mention d'une certaine molécule dans un article scientifique. Ils sont performants, même avec très peu de données d'entraînement.

Mais il y a aussi les Oracles d'apprentissage automatique ou les Oracles statistiques, qui peuvent être divisés en deux : les Oracles 'supervisés' et 'non supervisés'. Pour la création de modèles d'apprentissage automatique supervisés, les humains annotent les données d'entraînement avant de les envoyer à la machine. Chaque texte est jugé par au moins 3 humains: par exemple, s'il s'agit de spam ou non, s'il est positif ou négatif. Les Oracles d'apprentissage automatique non supervisés n'ont pas besoin de cette étape mais nécessitent de grandes quantités de données. C'est également à la machine de tracer ses propres motifs ou 'règles grammaticales'. Enfin, les experts font la différence entre les Oracles basés sur l'apprentissage automatique classique et ceux basés sur des réseaux de neurones. Vous en apprendrez plus à ce sujet dans la zone Lecteurs.

Les humains ont tendance à exagérer la performance des Oracles. Parfois, ces Oracles apparaissent quand il y a un dysfonctionnement. Dans les communiqués de presse, ces situations souvent dramatiques sont appelées des 'leçons'. Malgré la promesse de leurs performances, beaucoup de problèmes restent à résoudre. Comment s'assurer que les Oracles soient justes, que chaque être humain puisse les consulter, qu'ils soient compréhensibles par un large public ? Même au-delà, des questions existentielles persistent. Avons-nous besoin de tous les types d'intelligences artificielles ? Et qui définit ce qui est juste ou injuste ?

--- AdSense racial ---

Latanya Sweeney, professeur en Gouvernance et Technologie à l'Université de Harvard, a documenté une 'leçon' classique sur le développement des Oracles. En 2013, Sweeney, d'origine afro-américaine, a googlé son nom. Elle a immédiatement reçu une publicité pour un service qui lui offrait 'de voir le casier judiciaire de Latanya Sweeney'. Sweeney, qui n'a pas de casier judiciaire, a dès lors entamé une étude. Elle a commencé à comparer la publicité que Google AdSense offrait à différents noms racisés identifiables. Elle a découvert qu'elle recevait plus d'annonces de ce type en recherchant des noms ethniques non-blancs qu'avec des noms traditionnellement perçus comme blancs.

Sweeney a fondé son enquête sur des recherches portant sur 2184 prénoms racisés sur deux sites Web. 88 % des prénoms, identifiés comme étant donnés à un plus grand nombre de bébés noirs, sont considérés comme prédictifs de la race, contre 96 % de blancs. Les prénoms qui sont principalement donnés à des bébés noirs, comme DeShawn, Darnell et Jermaine, ont généré des annonces mentionnant une arrestation dans 81 à 86 % des recherches de noms sur un site, et dans 92 à 95 % des cas sur l'autre. Les noms qui sont principalement attribués aux blancs, comme Geoffrey, Jill et Emma, n'ont pas donné les mêmes résultats. Le mot 'arrestation' n'est apparu que dans 23 à 29 % des recherches de noms blancs sur un site, et 0 à 60 % sur l'autre.

Sur le site affichant le plus de publicité, un nom d'identification noir était 25 % plus susceptible d'obtenir une publicité suggérant un dossier d'arrestation. Quelques noms n'ont pas suivi ces modèles : Dustin, un nom donné principalement aux bébés blancs, a généré une publicité suggérant une arrestation dans 81 et 100 % des cas. Il est important de garder à l'esprit que l'apparition de l'annonce est liée au nom lui-même et non au fait qu'il ait un dossier d'arrestation dans la base de données de l'entreprise.

Référence : <https://dataprivacylab.org/projects/onlineads/1071-1.pdf>

--- Qu'est-ce qu'un bon employé ? ---

Depuis 2015, Amazon compte environ 575 000 travailleurs, et ils leur en faut plus. Par conséquent, ils ont mis sur pied une équipe de 12 personnes pour créer un modèle qui trouverait de bons candidats en parcourant des sites de demande d'emploi. L'outil attribuerait aux candidats une note allant de une à cinq étoiles. Le potentiel a alimenté le mythe : l'équipe voulait un logiciel qui recracherait les cinq meilleurs sur une liste de 100 candidats humains pour les embaucher. !!!

Le groupe a créé 500 modèles algorithmiques, centrés sur des fonctions et des lieux de travail spécifiques. Ils ont appris à reconnaître 50 000 termes qui figuraient sur les lettres d'anciens candidats. Les algorithmes ont appris à accorder peu d'importance aux compétences communes aux candidats en IT, comme la capacité d'écrire du code informatique, mais ils ont aussi reproduit les erreurs de leurs créateurs. Juste avant d'approuver un modèle, l'entreprise s'est rendue compte que les modèles ont décidé que les candidats masculins étaient préférables. Ils pénalisaient les candidatures qui comprenaient le mot 'femmes' ou 'féminin', comme dans 'capitaine de club d'échecs féminin'.

Et ils ont rétrogradé les diplômées de deux universités réservées aux femmes.

Ceci est dû à l'utilisation pour leur entraînement des demandes d'emploi reçues par Amazon sur une période de 10 ans. Durant cette période, l'entreprise avait surtout embauché des hommes. Au lieu de fournir la prise de décision 'équitable' que l'équipe d'Amazon avait promise, les modèles reflétaient une tendance biaisée dans l'industrie technologique. Mais ils l'ont aussi amplifiée et rendu invisible. Les activistes et les critiques affirment qu'il pourrait être extrêmement difficile de poursuivre un employeur en cas d'embauche automatisée : les candidats à un emploi pourraient ne jamais savoir que des logiciels intelligents ont été utilisés dans ce processus.

Référence : <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scrap-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

--- Quantification de 100 ans de stéréotypes sexuels et ethniques ---

Dan Jurafsky est le co-auteur de 'Speech and Language Processing', un des ouvrages les plus influents pour l'étude du traitement du langage naturel. Avec quelques collègues de l'Université de Stanford, il a découvert en 2017 que les 'word embeddings' peuvent être un outil puissant pour quantifier systématiquement les stéréotypes communs ainsi que d'autres tendances historiques.

Les 'word embeddings' sont une technique qui traduit les mots en vecteurs numérotés dans un espace multidimensionnel. Les vecteurs qui apparaissent proches l'un de l'autre, indiquent une signification similaire. Ainsi, tous les numéros seront regroupés, toutes les prépositions, les prénoms et les professions, etc. Cela permet de faire des calculs avec les mots. Vous pourriez, par exemple, soustraire Londres de Royaume-Unis et votre résultat serait le même que de soustraire Paris de France.

Un exemple de leur recherche montre que le vecteur de l'adjectif 'honorable' est plus proche du vecteur 'homme', alors que le vecteur 'soumis' est plus proche de 'femme'. Ces stéréotypes sont alors automatiquement appris par l'algorithme. Il s'avère problématique lorsque les 'embeddings' pré-entraînés sont utilisés pour des applications sensibles comme les classements de recherche, les recommandations de produits ou les traductions. Ce risque est réel, car un grand nombre de 'word embeddings' pré-entraînés sont téléchargeables sous forme de paquets prêts à l'emploi.

On sait que la langue reflète et maintient en vie les stéréotypes culturels. L'utilisation des 'word embeddings' pour repérer ces stéréotypes est moins cher et prends moins de temps que les méthodes manuelles. Mais leur mise en oeuvre dans des modèles de prédiction suscite beaucoup de discussions au

sein de la communauté du machine learning. Ces modèles fallacieux ou biaisés sont synonymes d'une discrimination automatisée. La question se pose : est-il vraiment possible d'éliminer complètement les préjugés de ces modèles ?

Certains affirment que oui, d'autres sont en désaccord. Avant de soumettre le modèle à une ingénierie inversée, nous devrions nous demander si nous en avons besoin tout court. Ces chercheurs ont suivi une troisième voie. En reconnaissant la discrimination qui trouve son origine dans le langage, ces modèles deviennent pour eux des outils de sensibilisation, en visualisant le problème.

L'équipe de la Stanford University a développé un modèle d'analyse des 'word embeddings' entraîné sur 100 ans de textes. Pour l'analyse contemporaine, ils ont utilisé les Google News word2vec Vectors, un paquet prêt à l'emploi, téléchargeable, entraîné sur le Google News Dataset. Pour l'analyse historique, ils ont utilisé des 'word embeddings' qui ont été entraînés sur Google Books et The Corpus of Historical American English (COHA <https://corpus.byu.edu/coha/>) avec plus de 400 millions de mots de textes des années 1810 à 2000. Afin de valider le modèle, ils ont entraîné des 'word embeddings' du New York Times Annotated Corpus pour chaque année entre 1988 et 2005.

Leur recherche montre que les 'word embeddings' reflètent l'évolution des stéréotypes sexistes et ethniques au fil du temps. Ils quantifient comment des préjugés spécifiques diminuent avec le temps tandis que d'autres stéréotypes augmentent. Les principales transitions révèlent des changements dans les descriptions de genre et de groupes ethniques lors du mouvement des femmes dans les années 1960-70 et la croissance de la population asio-américaine dans les années 1960 et 1980.

Quelques exemples :

Les dix professions les plus étroitement associées aux groupes ethniques dans le jeu de données de Google News :

- Hispanique : femme de ménage, maçon, artiste, concierge, danseur, mécanicien, photographe, boulanger, caissier, chauffeur.

- Asiatique : professeur, fonctionnaire, secrétaire, chef d'orchestre, physicien, scientifique, chimiste, tailleur, comptable, ingénieur.

- Blanc : forgeron, ferronnier, géomètre, shérif, tisserand, administrateur, maçon, statisticien, ecclésiaste, photographe.

Les 3 professions les plus masculines dans les années 1930 : ingénieur, avocat, architecte.

Les 3 professions les plus féminines dans les

années 1930 : infirmière, femme de ménage, aide-soignante.

Peu de choses ont changé dans les années 1990.

Principales professions masculines :
architecte, mathématicien et géomètre.
Les professions féminines restent les mêmes :
infirmière, femme de ménage et sage-femme.

Mais qu'est-ce qui s'est passé dans cette recherche avec les afro-américains?

Référence : <https://arxiv.org/abs/1711.08412>

--- Le Service ORES de Wikimedia ---

L'ingénieur de logiciels Amir Sarabadani a présenté le projet ORES à Bruxelles en novembre 2017 lors de notre Rencontre Algolittéraire. Cet 'Objective Revision Evaluation Service' utilise l'apprentissage automatique pour automatiser le travail critique sur Wikimedia, comme la détection du vandalisme et la suppression d'articles. Cristina Cochior et Femke Snelting l'ont interviewé.

Femke : Revenons à votre travail. Ces temps-ci, vous essayez de comprendre ce que signifie trouver des préjugés discriminatoires dans l'apprentissage automatique. La proposition de Nicolas Malevê, qui a donné l'atelier hier, était de ne pas essayer de le réparer, ni de refuser d'interagir avec des systèmes qui produisent de la discrimination, mais de travailler avec eux. Il considère que les préjugés sont inhérents à la connaissance humaine et que nous devons donc trouver des moyens de les utiliser d'une façon ou d'une autre. Nous avons discuté un peu de ce que cela signifierait, comment cela fonctionnerait... Je me demandais donc si vous aviez des idées sur cette question de partialité.

Amir : La partialité à l'intérieur de Wikipédia est une question délicate parce qu'elle se produit à plusieurs niveaux. Un niveau très discuté est le système des références. Toutes les références ne sont pas accessibles. Ce que la fondation Wikimedia a essayé de faire, c'est de donner un accès gratuit aux bibliothèques payantes. Ils réduisent l'exclusion en n'utilisant que des références en libre accès. Un autre type de discrimination est la connexion Internet, l'accès à Internet. Il y a beaucoup de gens qui ne l'ont pas. Une chose à propos de la Chine, c'est qu'Internet y est bloqué. Le contenu opposé au gouvernement de la Chine au sein du Wikipédia chinois est plus élevé parce que les éditeurs [qui peuvent accéder au site Web] ne sont pas pro-gouvernement et essaient de le rendre plus neutre. On le remarque donc à beaucoup d'endroits. En ce qui concerne l'intelligence artificielle (IA) et le modèle que nous utilisons chez Wikipedia, c'est plutôt une question de

transparence. Il existe un livre sur la façon dont les préjugés dans les modèles d'IA peuvent briser la vie des gens, intitulé 'Weapons of Math Destruction'. On y parle de modèles d'IA aux États-Unis qui classent les enseignants. C'est assez horrible parce qu'il y aura forcément des préjugés. D'après leur recherche, la façon d'aborder la question serait d'abord d'avoir un modèle open source, où l'on peut consulter le code et voir quelles fonctionnalités sont utilisées avec des données ouvertes, afin que les gens puissent enquêter, trouver des préjugés, donner leur feedback et faire un rapport. Il devrait y avoir un moyen de réparer le système. Je ne pense pas que toutes les entreprises vont dans cette direction, mais Wikipédia, en raison des valeurs qu'elle défend, est au moins plus transparente et pousse d'autres personnes à faire de même.

Référence : https://gitlab.constantvzw.org/algolit/algolit/blob/master/algoliterary_encounter/Interview%20with%20Amir/AS.aac

--- Tay ---

Une histoire tristement célèbre est celle du programme d'apprentissage automatique Tay, conçu par Microsoft. Tay était un chatbot qui imitait une adolescente sur Twitter. Elle a vécu moins de 24 heures avant d'être éteinte. Peu de gens savent qu'avant cet incident, Microsoft avait déjà entraîné et publié XiaoIce sur WeChat, l'application de chat la plus utilisée en Chine. Le succès de XiaoIce a été si prometteur qu'il a conduit au développement de son homologue américain. Cependant, les développeurs de Tay n'étaient pas préparés pour le climat de la plateforme Twitter. Bien que le bot savait distinguer un nom d'un adjectif, il n'avait aucune compréhension de la signification réelle des mots. Le robot a rapidement commencé à reproduire les insultes raciales et d'autres langages discriminatoires qu'il a appris par les autres utilisateurs de Twitter et les attaques de trolls.

L'apparition et la mort de Tay représentent une prise de conscience importante. Elle a montré les conséquences possibles de la corruption de l'apprentissage automatique, lorsque le contexte culturel dans lequel l'algorithme doit vivre n'est pas pris en compte.

Référence : <https://chatbotslife.com/the-accountability-of-ai-case-study-microsofts-tay-experiment-ad577015181f>

, i c i r e _i13%3 ++++++ 5 ++++++ '9 2 n 9s ea-silit e5 2- 14snn -c 8 c
ro3 s 3ro7 l l à sl D |c|l|e|a|n|e|r|s| |c|l|e|a|n| l 8 d768 88oe o+a 6 lail 71ea
aa 3 i V0tnt + u| ++++++ 8 ++++++ opti9 o 7 uu5ouc 7 1 8 r_'8 nd
a1 ft 8 ts a 7 tt n3ie i 6b ed+ r d 3a r u 9 rr2és0 p 23wv c 6st 2câ% i _7 a g |\n
a u3+ g+it | n é 1sr6 ot r 1rn6 t a o e - a' +/, t 9 i hl55ls4 t4 e r 2 t - a l t9
6el 4 4c n 79 xp -an_w2m+o r din o, +++++ n ++++++ ||6e g qes 6 2 s i m God76 e
oaed è h+ 4epe % p p ag lc w t 3 r |w|e| t s |h|e|l|p|e|d| _ 7r srt ed n % 9 9 eo19c ic ,
l3 4 u a pv i 9 u s- c ge7 +++++ r6 ++++++ r 7 or4 -s p 6e nl t x64 8)e t
+ 39 , - ru ê% o r c a +++++ cr ++++++ + - 4 w t 2r |s|l|n 7ad+ 95 D
pi t ' i e n, pu ai5h u i l t |w|e| u |c|l|e|a|n|e|d| a +4 d r cL ,4 s 9 5
éil o _i |i r n | i3ei e 1 +++++ \ y ++++++ d ii 4 b s 7 a e | re 0
t s u g98 4e 3, if uk ++++++ ++++++ ++++++ ++++++ r a s9 a 3| P 1 s9r
\e1 é i s is c |h|u|m|a|n| |w|o|r|k| - |i|s| |n|e|e|d|e|d| , r. ta pt 3 d0
r b d , an 9|an po d m ++++++ ++++++ at ++++++ ++++++ n -5 -+ o9 s Ds 1
i | d e471 ++++++ ++++++ 5 qel g8 2sM 4 4 | 8 - h jt5 8a 1 e
tr _l '| du w e5a |p|o|o|r|l|y|-|p|a|i|d| o l g ste a h2 s ai f 4 1 M 7
v 3a o6 s 15 ++++++ ++++++ s t 4 |s u / 1e t t Ae c re 90t46 r 8 eu
-ve r r-elr e ++++++ ++++++ | ++++++ ++++++ 5 r 8 4 d ,lcp s s
o i l 2r |f|r|e|e|l|a|n|c|e|r|s| u |c|a|r|r|y| |o|u|t| S m n 1: e8l o D
s |i t9 sm _ ++++++ ++++++ 5 ++++++ ++++++ 1 s 4 à 5 t d e
we lr a à ++++++ ++++++ s ++++++ ++++++ ++++++ h 4 5i
r r' dqe ean |v|o|l|u|n|t|e|e|r|s| L |d|o| |f|a|n|t|a|s|t|i|c| |w|o|r|k| u id t l
i o - s 9e 7 ++++++ ++++++ ++++++ ++++++ + u -
5 -i 9 4 t ++++++ ++++++ ++++++ ++++++ e 9 4 4 it
t l 7 r ré 2 |w|h|o|e|v|e|r| 1 |c|l|e|a|n|s| |u|p| |t|e|x|t| o n 6r a
sn l d n i esw ++++++ ++++++ ++++++ ++++++ 3 9 r p t n
ae i2 8 é _d 7 7 ê rs a c t e e 9r r6 i a
d il a t e n eê _ s e4t o L 6 ré s
i p nt 3c s é 8 2 e ot e a q l p cu e y
ot0e i 2 e a t c6 2 l m m u a
t 8 u - 6 h , i 4 u | 4 2 3 s +o o | i
r c t e7 n s 6 q a . u e a a 4
i 4s le i e 3 g a é u 8 9 - t c n 9 | - o , p
5 s r c d 3 t r s d i e t 8 o
1 . n p 5 l 2 b c8 r ét 3 c
a 4 8 w r s % c 3 e 9
5 s l 4 u 3 - w1 e s 3 st c +
. s nm 5 o 4 8
, h 2 s e a a g r tm 8 s gn
s c 3 oc e 1 a 2 u e t 2 cd
r 1 5 c lai % p4 e u 2
9 % / 6 34 u 9 a m s el9 uu
7 e s 4 e . e 9 a e e l 6 ô
e h o ng 2 e e ,
u ip e 4 p e h e o y o
c u 5 5 oa t5 i 5 e t 2
+ c u f l n . é u s 2 e t
w 6 a 5 é u s 2 e t
6 n l :| i o l a - w pa
8 + è i cl w i e / j 2
à u \ o s s D d 9 j
m d 9 p 7
t + 1 g 3 7s
' t d r 2
| p s 1 l e c 3
c q n v e
a 2 r
4 4 s e o é +m y l D a
sg a s - q e g s E
d a9 . g r' 6
6 i + s u r k e a s- V d r'
ea 3 d e r - r r g s a | 8 .

--- Projet Gutenberg et
Distributed Proofreaders ---

Le projet Gutenberg est notre grotte d'Ali Baba. Il offre plus de 58 000 livres électroniques gratuits à télécharger ou à lire en ligne. Les œuvres sont acceptées sur Gutenberg lorsque leur droit d'auteur américain a expiré. Des milliers de bénévoles numérisent et relisent des livres pour aider le projet. Une partie essentielle du travail est réalisée dans le cadre du projet Distributed Proofreaders. Il s'agit d'une interface Web pour aider à convertir les livres du domaine public en livres électroniques. Pensez aux fichiers texte, aux e-pubs, aux formats Kindle. En divisant la charge de travail en pages individuelles, de nombreux bénévoles peuvent travailler sur un livre en même temps, ce qui accélère le processus de nettoyage.

Pendant la relecture, les bénévoles reçoivent une image scannée de la page et une version du texte, lue par un algorithme de reconnaissance optique des caractères (OCR) entraîné pour reconnaître les lettres dans les scans. Cela permet de comparer facilement le texte à l'image, de le relire, de le corriger et de le renvoyer sur le site. Un deuxième bénévole se voit ensuite présenter le travail du premier. Il vérifie et corrige le travail si nécessaire, et le soumet au site. Le livre passe ensuite par un troisième cycle de relecture et deux autres cycles de mise en page à l'aide de la même interface Web. Une fois que toutes les pages ont terminé ces étapes, un post-processeur les assemble soigneusement dans un e-book et les soumet à l'archive du Projet Gutenberg.

Nous avons collaboré avec le Distributed Proofreaders Project pour nettoyer les fichiers numérisés que nous avons reçus de la collection du Munda-neum. De novembre 2018 jusqu'à la première mise en ligne du livre 'L'Afrique aux Noirs' en février 2019, An Mertens a échangé environ 50 courriels avec Linda Hamilton, Sharon Joiner et Susan Hanlon, toutes bénévoles du Distributed Proofreaders Project. La conversation complète est publiée en ligne. Cela pourrait vous inspirer à partager des livres non disponibles en ligne.

--- Une version algolittéraire
du Manifeste sur l'entretien ---

En 1969, un an après la naissance de son premier enfant, l'artiste new-yorkaise Mierle Laderman Ukeles a écrit un 'Manifesto for Maintenance' (Manifeste pour l'entretien).

Le Manifeste d'Ukeles appelle à une réévaluation de l'état des travaux d'entretien dans l'espace privé, domestique et public. Ce qui suit est une version modifiée de son texte inspirée par le travail des Nettoyeurs.

IDÉES

A. L'instinct de Mort et l'instinct de Vie :

L'Instinct de Mort : séparation ; catégorisation ; avant-garde par excellence ; suivre le chemin pré-dit vers la mort - exécuter son propre code ; changement dynamique.

L'Instinct de Vie : l'unification ; le retour éternel ; la perpétuation et l'ENTRETIEN de la matière ; les systèmes et opérations de survie ; l'équilibre.

B. Deux systèmes de base :

Développement et entretien. La boule de cristal de chaque révolution : après la révolution, qui va essayer de repérer le taux de discrimination dans la production ?

Développement : pure création individuelle ; le nouveau ; le changement ; le progrès ; l'avancée ; l'excitation ; la fuite ou s'enfuir.

Entretien : garder la poussière de la création individuelle pure ; préserver le nouveau ; soutenir le changement ; protéger le progrès ; défendre et prolonger l'avancée ; renouveler l'excitation ; répéter le vol ; montrez votre travail/remontez-le ; gardez le dépôt git mis à jour ; gardez l'analyse des données révélatrice.

Les systèmes de développement sont des systèmes de rétroaction partielle avec une grande marge de changement.

Les systèmes d'entretien sont des systèmes à rétroaction directe avec peu de possibilités de modification.

C. L'entretien est une corvée,
ça prend tout le temps.

L'esprit est éblouissant et s'irrite devant l'ennui.

La culture attribue un statut médiocre aux emplois d'entretien = salaire minimum, les Mechanical Turks d'Amazon = pratiquement aucun salaire.

Nettoyer le set, marquer les données d'entraînement, corriger les fautes de frappe, modifier les paramètres, terminer le rapport, satisfaire le demandeur, télécharger la nouvelle version, joindre les mots qui ont été mal reconnus par le logiciel de Reconnaissance Optique de Caractères, accomplir ces tâches d'intelligence humaine, essayez de deviner la signification du formatage du demandeur, vous devez accepter le 'hit' avant de pouvoir soumettre les résultats, résumer l'image, ajouter la

case de délimitation, quelle est la similitude sémantique de ce texte, vérifiez la qualité de la traduction, collecter vos micro-paiements, devenir un Mechanical Turk à succès.

Référence : <https://www.arnolfini.org.uk/blog/manifesto-for-maintenance-art-1969>

--- Une panique robotique chez le Mechanical Turk d'Amazon ---

Le Mechanical Turk d'Amazon prend le nom d'un automate d'échecs du 18ème siècle. En fait, le Turc mécanique n'était pas du tout une machine. C'était une illusion mécanique qui permettait à un maître d'échecs humain de se cacher à l'intérieur de la boîte et de l'utiliser manuellement.

Pendant près de 84 ans, le Turc a remporté la plupart des matchs joués lors de ses manifestations en Europe et en Amérique. Napoléon Bonaparte se serait lui aussi laissé berné par cette ruse.

Le Mechanical Turk d'Amazon est une plateforme en ligne à destination des humains pour exécuter des tâches que les algorithmes ne parviennent pas à faire. Il peut s'agir, par exemple, d'annoter des phrases comme étant positives ou négatives, de repérer des plaques d'immatriculation, de reconnaître des visages. Les postes affichés sur cette plateforme sont souvent rémunérés moins d'un centime par tâche. Les tâches les plus complexes ou nécessitant le plus de connaissances peuvent être payées jusqu'à plusieurs centimes. Pour gagner leur vie, les 'turkers' doivent accomplir le plus de tâches possible le plus rapidement possible, ce qui entraîne d'inévitables erreurs. Les créateurs des jeux de données doivent incorporer des contrôles de qualité lorsqu'ils publient un travail sur la plate-forme. Ils doivent vérifier si le 'turker' a réellement la capacité d'accomplir la tâche, et ils doivent également vérifier les résultats. De nombreux chercheurs universitaires utilisent le Mechanical Turk pour des tâches qui auraient été exécutées par des étudiants auparavant.

En août de l'année dernière, Max Hui Bai, un étudiant en psychologie de l'Université du Minnesota, a découvert que les enquêtes qu'il a menées avec Mechanical Turk étaient pleines de réponses absurdes aux questions ouvertes. Il a retracé les mauvaises réponses et a découvert qu'elles avaient été soumises par des répondants ayant des coordonnées GPS en double. Cela a suscité des soupçons.

Bien qu'Amazon interdise explicitement aux robots d'effectuer des travaux sur Mechanical Turk, l'entreprise ne publie pas les problèmes qu'ils causent sur sa plate-forme. Les forums pour 'turkers' sont pleins de conversations sur l'automatisation du travail, le partage de pratiques sur la façon de créer des robots qui transgresseraient les termes d'Amazon. Vous pouvez également trouver

des vidéos sur YouTube montrant aux 'turkers' comment écrire un bot qui remplit des réponses pour vous.

Kristy Milland, une militante de Mechanical Turk, dit : 'Les travailleurs sur Mechanical Turk ont été très, très mal traités pendant 12 ans et, d'une certaine façon, je vois cela comme un point de résistance. Si nous étions payés équitablement sur la plateforme, personne ne prendrait le risque de perdre son compte de cette façon.'

Bai a créé un questionnaire pour les chercheurs en dehors de Mechanical Turk. Il dirige actuellement une recherche parmi les spécialistes des sciences sociales pour déterminer la quantité de données erronées utilisées, l'ampleur du problème et les moyens de l'enrayer. Mais il est impossible à l'heure actuelle d'estimer combien de jeux de données sont devenus peu fiables de cette façon-ci.

Références :

<https://www.wired.com/story/amazon-mechanical-turk-bot-panic/>

<https://www.maxhuibai.com/blog/evidence-that-responses-from-repeating-gps-are-random>

<http://timryan.web.unc.edu/2018/08/12/data-contamination-on-mturk/>

ile a3ra 4rmo 24c ++++++ s ++++++ ny. 2hn si l 7sé n9r-
pc- , 92é né e r | |i|n|f|o|r|m|a|n|t|s| |i|n|f|o|r|m| 5e 32| a m l+m r r
em 'r i,e e c én ++++++ , ++++++ 3 iso fIc éT o 8 ne3é3e-
t mm i 4e G 6 o i s 3 4 u i t 3r -r é 8 m r ri r _ di453 f d c e u i -_tm- o | e n+ 9e
L u 8 S n 9 r + s e u 2 3 5 C e m e 8 u , é 9 r 2 9 u s l u 2 e s m ê r r o l 6 u r d e o s 3+ e r
j . e 5 1 i é d ++++++ r ++++++ ++++++ 7 p p i a e
m o i 6 9 6 r 5 o l h 5 0 |e|a|c|h| |d|a|t|a|s|e|t| |c|o|l|l|e|c|t|s| |d|i|f|f|e|r|e|n|t| p l o e
1 l i % m 1 - i m t ++++++ t ++++++ ++++++ s l b a i s / _ t
P ` n c r 3 w d _ t q % i c _ c i e 3 a n - i ++++++ ++++++ _ 1 î t C t o i n e
i p r m l l _ + 4 n n 2 f i n e é d c \ _ |i|n|f|o|r|m|a|t|i|o|n| |a|b|o|u|t| o l t r m s 8 n
e \ o m | - o i e o % a 8 é 6 t t i 5 s v ++++++ ++++++ t e r i i s e
r i 4 l 2 sé n - t h n 1 m 6 t l - d 5 4 _ l s b ++++++ ++++++ a o r i q 5 l z 8
o h e w + a 6 1 t o o 3 | b l e u |t|h|e| |w|o|r|l|d| i p y e s g j _
d a i - - t s e f i s n 1 4 n l + r 0 r - ++++++ ++++++ m r _ r l D r o t s
3 l n i s l q e a r u i ++++++ a l r ++++++ ++++++ ++++++ i _ a 2
w o 7 3 6 n , s c i - |d|a|t|a|s|e|t|s| |a|r|e| |i|m|b|u|e|d| |w|i|t|h| s e e a e i 3
i 6 h % 2 l p 8 w d 5 l ++++++ ++++++ ++++++ e x r s e u
. é r g t a n o 3 f n 5 u ++++++ ++++++ ++++++ _ r 2 a 8 s o
d n + n r 9 a 1 e e i 5 d |c|o|l|l|e|c|t|o|r|'|s| |b|i|a|s| u 8 8 6 u o
n t a e s i s i m n 4 2 c e % t t 7 6 l 4 r é ++++++ ++++++ | d e i m t t 3 7 r
s e 8 e e ++++++ ++++++ ++++++ o r (++++++ ++++++ ++++++ s 9 t
s h s _ 9 1 L |s|o|m|e| |d|a|t|a|s|e|t|s| |c|o|m|b|i|n|e| |m|a|c|h|i|n|i|c| 7 4 v 7
m n c p | ? ++++++ ++++++ ++++++ u ++++++ ++++++ ++++++ 1 â r 3 u d p
i v c 3 n i a c i é é e ' e r 5 + | c / a a q ++++++ ++++++ ++++++ + o 7 r
% r n o + e r u o w i 3 n |l|o|g|i|c| |w|i|t|h| |h|u|m|a|n| a l s
a u r , o t s - 7 r v t % s e ++++++ ++++++ ++++++ e 0 m n
W r 4 s - e 0 1 2 a r 6 n t | ++++++ , s s e r 1 | i
p f i 1 u 5 i i a r o + s e |l|o|g|i|c| 5 l s c p d \
e e e l e q a s e s r n c ++++++ | t _ \ o
e i e % c ++++++ ++++++ ++++++ e w 7 8
i , _ j t l |m|o|d|e|l|s| |t|h|a|t| s p o a i S é r é u e , 9 s
5 - r 6 p r o ++++++ ++++++ ++++++ o 7 9 9 9 7 / | u s q
o 1 0 d s i ++++++ ++++++ ++++++ i 3 e e a ê p i m e t
r h é e m o 1 |r|e|q|u|i|r|e| c 9 r 1 e t p t
e p c % o 8 h ++++++ ++++++ ' f 4 é 5 u a 9 s B
- c s i 5 ++++++ ++++++ ++++++ r ++++++ ++++++ ++++++ e
r u 3 t s t l |s|u|p|e|r|v|i|s|i|o|n| |p |m|u|l|t|i|p|l|i|y| |t|h|e| b l e s u 5 p r
h | ++++++ ++++++ ++++++ ++++++ ++++++ ++++++ t u e n o
o e 5 m a - e u h 6 ++++++ ++++++ ++++++ i e q
n c c r \ n 9 |s|u|b|j|e|c|t|i|v|i|t|i|e|s| 8 _ e 7 s
c U f s n ô e ++++++ ++++++ ++++++ 3 5 i t |
e e g r r s s d ++++++ ++++++ ++++++ - l a o
e e + m r 4 8 |m|o|d|e|l|s| c |p|r|o|p|l|a|g|a|t|e| |w|h|a|t|
. e - \ o ++++++ ++++++ ++++++ ++++++ ++++++ m l 6 t
1 p v l _ s r ++++++ ++++++ ++++++ V
m i u 1 t t a v |t|h|e|y|' |v|e| |b|e|e|n| r 3 e ' c e
9 e s e 8 o t g g t ++++++ ++++++ ++++++ c f 4 r q s
7 _ - o 3 1 9 d ++++++ ++++++ ++++++ r t e
2 n e 7 n g ++++++ ++++++ ++++++ - 0 q u p
1 u ++++++ ++++++ ++++++ r e x u 7 é
9 a ' é |s|o|m|e| |o|f| |t|h|e| 2 a 1 1 e
e u p ++++++ ++++++ ++++++ t 1 1 e
2 s e ++++++ ++++++ ++++++ ++++++ ++++++ ++++++
u _ |d|a|t|a|s|e|t|s| |p|a|s|s| |a|s| |d|e|f|a|u|l|t| |i|n|
l 4 f p i ++++++ ++++++ ++++++ ++++++ ++++++ ++++++ t
o c e d 6 t ++++++ ++++++ ++++++ 1 5
2 pi |t|h|e| |m|a|c|h|i|n|e| 6 e d
I , 7 p ++++++ ++++++ ++++++ l A
- 4 n s ++++++ ++++++ ++++++ ++++++ a
r é m e p o r _ t D |l|e|a|r|n|i|n|g| |f|i|e|l|d| n n é
s , ++++++ ++++++ ++++++ p ++++++ ++++++ ++++++ h
8 t d |h|u|m|a|n|s| |g|u|i|d|e| |m|a|c|h|i|n|e|s| c t
a m o 1 ++++++ ++++++ ++++++ ++++++ ++++++ i
r i . g | c Q b 7
o 3 8 é o | 8 n p 7 9 e o è t e
a r 8 e 8 n 7 9 e o o a

--- Les jeux de données comme représentations ---

Les processus de collecte des données qui mènent à la création du jeu de données soulèvent des questions importantes : qui est l'auteur des données ? Qui a le privilège de collectionner ? Pour quelle raison la sélection a-t-elle été faite ? Que manque-t-il ?

L'artiste Mimi Onuoha donne un exemple excellent de l'importance des stratégies de collection. Elle choisit le cas des statistiques relatives aux crimes haineux. En 2012, le Programme de déclaration uniforme de la criminalité (DUC) du FBI a enregistré 5 796 crimes haineux. Toutefois, le Bureau a établi 293 800 rapports sur de tels cas. C'est plus de 50 fois plus. La différence entre les chiffres peut s'expliquer par la façon dont les données ont été recueillies. Dans le premier cas, les organismes d'application de la loi de tout le pays ont volontairement signalé des cas. Pour le deuxième, le Bureau des statistiques a distribué l'enquête nationale sur la victimisation directement aux foyers des victimes de crimes motivés par la haine.

Dans le domaine du traitement du langage naturel, le matériel avec lequel les modèles d'apprentissage automatique travaillent est le texte, mais les mêmes questions se posent : qui sont les auteurs des textes qui composent les jeux de données ? Au cours de quelle période les données ont-elles été recueillies ? Quel type de vision du monde représentent-elles ?

En 2017, l'algorithme Top Stories de Google a placé un fil de discussion trompeur du site 4chan en haut de la page de résultats lors de la recherche du tireur de Las Vegas. Le nom et le portrait d'une personne innocente étaient liés au crime. Bien que Google ait changé son algorithme quelques heures seulement après que l'erreur ait été découverte, cela a sérieusement affecté la personne. Une autre question persiste : pourquoi Google n'a-t-il pas exclu le site de ragôts 4chan du jeu des données d'entraînement ?

Références :

<https://points.datasociety.net/the-point-of-collection-8ee44ad7c2fa>

<https://arstechnica.com/information-technology/2017/10/google-admits-citing-4chan-to-spread-fake-vegas-shooter-news/>

--- L'annotation pour un Oracle qui détecte le vandalisme sur Wikipédia ---

Ce fragment est extrait d'une interview avec Amir Sarabadani, ingénieur de logiciels chez Wikimedia.

Il était à Bruxelles en novembre 2017 lors de la Rencontre Algolittéraire.

Femke : En considérant Wikipedia comme une communauté vivante, chaque nouvelle page change le projet. Chaque modification est en quelque sorte une contribution à un organisme vivant de la connaissance. Donc, si au sein de cette communauté vous essayez de distinguer ce qui rend service à la communauté et de généraliser ceci dans un modèle - car je pense que c'est ce que l'algorithme de la bonne ou mauvaise foi essaie de faire - vous le faites sur base d'une généralisation de l'idée abstraite de Wikipedia, et non sur base de l'organisme vivant. Ce qui m'intéresse dans la relation entre le vandalisme et ce débat, c'est la façon dont nous pouvons comprendre la dynamique conventionnelle de ces processus d'apprentissage automatique. Si on distingue la bonne ou la mauvaise foi sur base d'étiquettes préexistantes et qu'on la reproduit ensuite dans des modèles algorithmiques, comment tenir compte des changements qui se produisent, c'est-à-dire de la vie réelle du projet?

Amir : C'est une discussion intéressante. Premièrement, ce que nous appelons la bonne ou la mauvaise foi provient de la communauté elle-même; nous ne faisons pas l'annotation nous-mêmes, c'est la communauté qui le fait. Ainsi, dans beaucoup de Wikipédias de langues différentes, la définition de ce qui est la bonne ou la mauvaise foi sera différente. Wikimedia essaie de refléter ce qui se trouve à l'intérieur de l'organisme et non de changer l'organisme lui-même. Si l'organisme change et que nous constatons que la définition de la bonne foi à Wikipédia a été modifiée, nous mettons en œuvre cette boucle de rétroaction qui permet aux gens de porter un jugement sur leurs modifications à l'intérieur de leur communauté. S'ils sont en désaccord avec l'annotation, nous pouvons revenir au modèle et modifier l'algorithme pour refléter ce changement. C'est une sorte de boucle fermée : vous changez les choses et si quelqu'un voit qu'il y a un problème, il nous le dit et nous pouvons modifier l'algorithme. C'est un projet en cours.

Référence : https://gitlab.constantvzw.org/algolit/algolit/blob/master/algoliterary_encounter/Interview%20with%20Amir/AS.aac

--- Comment faire connaître votre jeu de données ---

NLTK signifie Natural Language Toolkit. Pour les programmeurs qui traitent le langage naturel avec Python, c'est une bibliothèque essentielle. De nombreux rédacteurs de tutoriels recommandent aux programmeurs d'apprentissage automatique de commencer par les jeux de données NLTK intégrés. Il compte 71 collections différentes, avec un total de près de 6000 éléments.

Parmi eux, on trouve le corpus Movie Review pour l'analyse des sentiments. Ou le corpus Brown, qui a été créé dans les années 1960 par Henry Kučera et W. Nelson Francis à l'Université Brown de Rhode Island. Il y a aussi le corpus de la Déclaration des droits de l'homme, qui est couramment utilisé pour vérifier si un code peut fonctionner dans plusieurs langues. Le corpus contient la Déclaration des droits de l'homme dans 372 langues du monde entier.

Mais quel est le processus pour faire accepter un jeu de données dans la bibliothèque NLTK de nos jours ? Sur la page Github, l'équipe nltk décrit les exigences suivantes :

- Ne rajoutez que les corpus qui ont obtenu un niveau de notabilité de base. Cela signifie qu'il existe une publication qui le décrit et une communauté de programmeurs qui l'utilisent.

- Assurez-vous d'avoir l'autorisation de redistribuer les données et de pouvoir les documenter. Cela signifie qu'il est préférable de publier le jeu de données sur un site Web externe avec une licence.

- Utilisez les lecteurs de corpus NLTK existants lorsque c'est possible, ou bien apportez un lecteur de corpus bien documenté à NLTK. Cela signifie que vous devez organiser vos données de manière à ce qu'elles puissent être facilement lues à l'aide du code NLTK.

Référence : <http://www.nltk.org/>

--- Extrait d'une critique positive d'un film IMDb du jeu de données NLTK ---

corpus : movie_reviews

fichier : pos/cv998_14111.txt

le deuxième film épique de steven spielberg sur la seconde guerre mondiale est un chef-d'œuvre incontesté du cinéma . spielberg , encore étudiant en cinéma , a réussi à ressusciter le genre de la guerre en produisant l'un de ses films les plus poignants et les plus puissants . il a également réussi à faire briller tom hanks , qui livre une performance époustouflante . pendant environ 160 de ses 170 minutes, ' sauver le soldat ryan ' est sans faille . littéralement . l ' histoire est assez simple . après l ' invasion du jour J (dont les séquences sont tout à fait spectaculaires), capt . john miller (joué par tom hanks) et son équipe sont forcés à chercher un soldat . james ryan (joué par matt damon), dont les frères sont tous morts au combat. une fois qu ' ils l ' ont trouvé , ils doivent le ramener immédiatement pour qu'il puisse rentrer chez lui . la compagnie de miller est composée d ' acteurs aux jeux tout sim-

plement sensationnels : bary pepper , adam goldberg , vin diesel , giovanni ribisi , davies et burns . le film se clôture avec des scènes de bataille extraordinaires .

--- Les ouroboros de l'apprentissage automatique ---

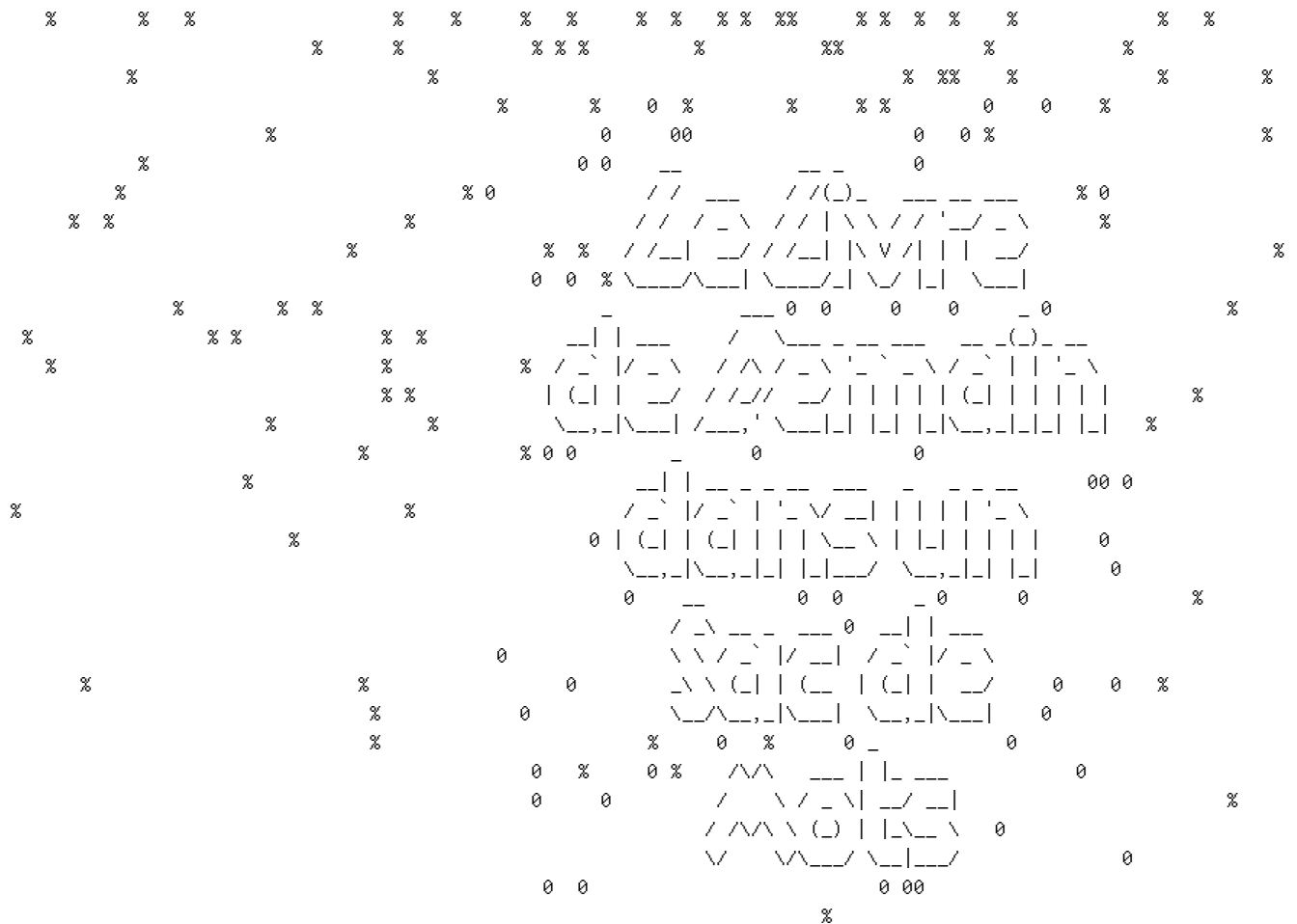
Wikipédia est devenue une source d'apprentissage non seulement pour les humains, mais aussi pour les machines. Ses articles sont des sources de premier ordre pour l'entraînement de modèles. Le matériel avec lequel les machines sont entraînées est identique au contenu qu'elles ont aidé à écrire. En fait, au début de Wikipédia, de nombreux articles ont été écrits par des robots. Rambat, par exemple, était un robot controversé sur la plateforme anglophone. Il est l'auteur de 98% des pages décrivant les villes américaines.

A cause de ces interventions de robots thématiques et régulières, les modèles de prédiction qui sont entraînés sur le dump de Wikipedia ont une vision unique de la composition des articles. Par exemple, un modèle thématique entraîné sur l'ensemble des articles de Wikipédia associe 'rivière' à 'Roumanie' et 'village' à 'Turquie'. C'est parce qu'il y a plus de 10000 pages écrites sur les villages en Turquie. Cela devrait suffire à susciter des envies de voyage, mais c'est bien trop par rapport à d'autres pays. L'asymétrie provoque une fausse corrélation et doit être corrigée. La plupart des modèles tentent d'exclure le travail de ces auteurs robots prolifiques.

Référence : <https://blog.lateral.io/2015/06/the-unknown-perils-of-mining-wikipedia/>

0 12 3 4 5 67 8 9 0
 12 3 4 5 67 8 9 0 12
 3 4 5 67 8 9 0 1 2 3
 4 56 7 8 9 01 2 3
 4 56 7 8 9 01 2 3 4
 5 6 7 8 9 0 1 2 3 4 5 6
 7 8 9 0 1 2 3 4 5 6 7 8 9
 7 89 0 1 2 34 5 6 7 89
 89 0 1 2 3 4 5 6 7 8 9
 0 1 23 4 5 6 78 9 0
 1 2 3 4 5 6 78 9 0
 1 2 3 4 5 6 7 8 9 0 12
 3 4 5 67 8 9 0 12 3
 4 5 6 7 8 9 0 1 2 3
 4 56 7 8 9 01 2 3 4
 5 6 7 8 9 0 1 2 3 4 5 6
 7 8 9 0 1 2 3 4 5 6 7 8 9
 7 8 9 0 1 2 3 4 5 6 7 89
 89 0 1 2 34 5 6 7 89
 0 1 2 3 4 5 6 7 8 9 0
 1 2 3 4 5 6 7 8 9 0 12 3
 4 5 6 7 8 9 0 1 2 3
 4 5 6 7 8 9 0 1 2 3 4
 56 7 8 9 01 2 3 4 5
 6 7 8 9 0 1 2 3 4 5 6
 7 8 9 0 1 2 3 4 5 6 7 89
 7 8 9 0 1 2 3 4 5 6 7 89
 0 1 2 34 5 6 7 89 0
 1 2 34 5 6 7 8 9 0
 1 23 4 5 6 78 9 0
 1 23 4 5 6 7 8 9 0 12 3
 2 3 4 5 6 7 8 9 0 12 3
 4 5 6 7 8 9 0 12 3

9 nl i 5 ' r c ++++++ a ++++++ -o on an r e c coeun b 9 em t
-rm6 n r r 5ui rt s r ar |r|e|a|d|e|r|s| mg |r|e|a|d| s-dn + 1 tm n u7| 1e0+ iuae3Mii+u e3 l
i e t e t % c_ | 5l ++++++ qc ++++++ nem e ic-7- ro + g i_ 6é
1 4i + c s r74 ss a_ 9é- S t 8ra | 3 t l nr , mnr 6 i l 8 + 4'99 m u a t7% s_
é 2 8 un e 8/n 2m c5 2nsr u 0 2t) 6 v 5 lu6je_r a \ r su r5 9 sdaur o eu h p
r 3 e tes 7ei I t% 9 +++ ++++++ | ++++++ 6 | /s l3 e pu5 n e
47 Mr t rd l ' + n 3 |a| |c|o|m|p|u|t|e|r| u |u|n|d|e|r|s|t|a|n|d|s| +- c 7r s ph r o6 9a
do 1 2² 1 | +++ ++++++ 5 ++++++ e n s nse tk l â0 ' é|
d é. s ees e n ++++++ tLo ++++++ Cnr1 5l _ rUt g od
17 s8 | s 9 8+ 7 |a|l|l| |m|o|d|e|l|s| u |t|r|a|n|s|l|a|t|e| % s_ a n t 2 , - r rer
uia1 p et _ |n s 9_l1 u ++++++ % ++++++ w ae m ssI |o i um.i é|r
sarpé u m ou_ | , o ++++++ r l ++++++ 1hff 0 g M l a 5,7 |y
s Sr c r xâd 5 e |s|o|m|e| |m|o|d|e|l|s| f |c|o|u|n|t| 3nrpds6 eo a t i rn e c + 7 n
1 V d e i e e -5;d ++++++ è ++++++ l x n é r45er é -a i _ l 2 i o 61
'h e| et - 9 e\ +s ++++++ We ++++++ t s t s 5 2 i 1 | etu
+ld t er |s|o|m|e| |m|o|d|e|l|s| |r|e|p|l|a|c|e| s p e. _ r8ai 7 ss n - g s
nc (- 1 e u. ++++++ e h i \ Sr_e ie _ lés
aa% rns1 .7,a m 5 a p n 6 s 5 n e s 8 8 | 4 77 n i ttl a ,
6 t e o le r- oa e l3i _ n e ws lo i i C a t66r r ses n2 5q 6 s e
7 e nn s l 5 e s_ év r o e t +\p \ 7 2 it e 1 e 8 t %
. a) e a ee8 emu r | eun t n s 8eel e o 9 p s eei e ee l
r 1 '6 ts teae 4 a , n n t 66 n o n e3o e èla en tu + , hw 8 18
6 5é e t i w rm x+ 2 ç e td eu aen % +i nr 76 c 59u 4u ea
u 7 cse ise i oa m _ no) t g h 3 2 4 ll u 7l
à ad l a 64,F As9 m% s s f w8 ts t 3 s7m t a ndh h4 3 n +s t n
4 |s6 6 1rr r ne d D r_1 - 33 2 \o 2 o r m5 s _5n 9ii a
8i 4 465 _ i sé 1 2 s f r9 l e o j rii e le 6 c t ep
V u t 7le s / 2 d 89 t 6r 2 , b r ht r p , 4lt e
, 1 3d s é - t 2 e u e oe . u u. s 7 e 3
rdt i 4 n 4 8 n m p o l a 8 r 7 i w
n s V1 he ud T M| o a A n 4 5 u 5 c 0 ua h4 s e de
_7 c 3 a h n 8 5 u u , n ê p d n v u 9
r i e i i 4 f | 7 c d d 8 s r + r
s m 8 v i 0 s , r - 5 5 r 1 e t
' e d 3 6 nz - e et na 7 | n n s i aa s e s e
p t , +d r 0 2 a be / n id a gu % + c ' id 8e r m
s sv - c r _ ;u 7 i_ 6 % i s é a d e
a t tn i + l a r s à d 2 o 1 b _ t 7t e nn c
7 en e a 1 s u é h 7q e i 39 me c 1 m
s é 0 9 u 4 l l e r 62 vs tk ne 6 ee
d a s 2 ne N r à | 4 9 s
de a a 4 1 3 o 7 r 4 9
6 , r u a t ' n m f es e c / e
m p 1 a8p k e , n e t _ 9 r t \
r 8 u n _ 5 - 8 et p m V 2 r 3 e
. t + e ir \ c e , 5 pN
p l lp | _ as l e o 8 p
o b tu 8 m d 5 93 s8 c c s 4 .
s 6 5 9 0 et . d i
t p% 7 -i n e 1
4 r r m 0 l e , e u f a e4 v f
7 é 6 9 6é | \ é + \ 0u - e a
a t u 3a , , - b
| uu - I e o 4 3 r r x à d 1 9 o
i V r 4 r ê 9 i u
o a o + ' c i 8 6 o
p i 9 6 % e s % c sn 8 e s
i l 2g i n c / m a t se n l
n h i s i a _ r 93 + e
1 h i s + p 4 S d6 c
r és 3 3 % n)
l | i l d a R
5m h i 6 n o 6
m s p 6 e - t i



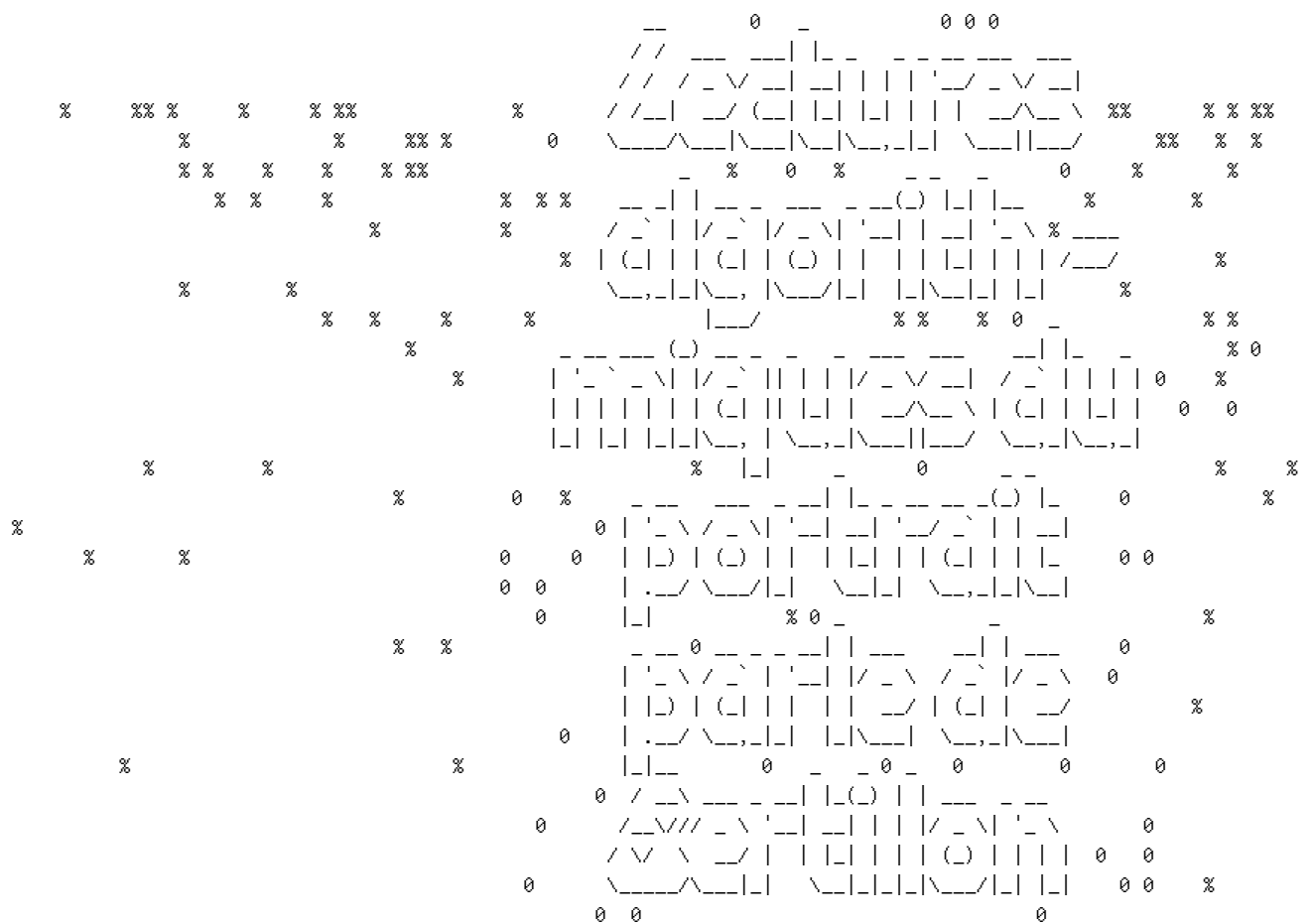
% par Algolit

%
 Le modèle du 'sac de mots' est une représentation simplifiée du texte utilisé dans le traitement du langage naturel. Dans ce modèle, un texte est représenté sous forme de collection de mots uniques, sans tenir compte de la grammaire, de la ponctuation et même de l'ordre des mots. Le modèle transforme le texte en une liste de mots et leur occurrence dans le texte, littéralement un sac de mots.

%
 Cette forte réduction de la langue fut un choc au début de nos expériences en apprentissage automatique. Le sac de mots est souvent utilisé comme référent, sur base duquel le nouveau modèle doit s'efforcer d'être plus performant. Il peut comprendre le sujet d'un texte en reconnaissant les mots les plus fréquents ou importants. On mesure souvent les similitudes des textes en comparant leurs sacs de mots. %

% Pour cet ouvrage, l'article 'Le Livre de Demain' de l'ingénieur G. Vander Haeghen, publié en 1907 dans le Bulletin de l'Institut International de Bibliographie, a été littéralement réduit à un sac de mots. Vous pouvez acheter votre exemplaire à l'accueil du Mundaneum.

Concept & réalisation: An Mertens



par Guillaume Slizewicz (Espèces urbaines)

'Un code télégraphique du portrait parlé', écrit en 1907, est une tentative de traduire en chiffres le 'portrait parlé', technique de description du visage créée par Alphonse Bertillon, créateur de l'anthropométrie judiciaire. En appliquant ce code, Otlet espérait que les visages des criminels et des fugitifs pourraient être facilement communiqués par voie télégraphique. Dans sa forme, son contenu et son ambition, ce texte représente la relation complexe que nous entretenons avec les technologies documentaires. Ce document a été choisi comme base pour la création des installations suivantes pour trois raisons.

- Premièrement, ce texte est un algorithme en soi, un algorithme de compression, ou pour être plus précis, la présentation d'un algorithme de compression. Il tente de réduire la taille de l'information tout en la gardant lisible pour la personne possédant le code. À cet égard, elle est étroitement liée à la façon dont nous créons notre technologie, à la recherche d'une plus grande efficacité, de résultats plus rapides et de méthodes moins coûteuses. Il représente notre appétit de chiffrement qui s'étend au monde entier, notre envie de mesurer les plus petites choses, d'étiqueter les différences les plus infimes... Ce texte incarne en lui-même la vision du Mundaneum.

- Deuxièmement, on y traite des raisons et des mises en œuvre de nos technologies. La présence de ce texte dans les archives sélectionnées est presque ironique à une époque où la reconnaissance faciale et la surveillance des données font la une des journaux. Ce texte présente les mêmes caractéristiques que certaines technologies d'aujourd'hui : il est motivé par un contrôle social, classifie les personnes, pose les bases d'une société de surveillance. Les caractéristiques physiologiques sont au cœur de récentes controverses : les photos d'identité ont été standardisées par Bertillon, elles sont maintenant utilisées pour entraîner des réseaux neuronaux à identifier les criminels, les systèmes

Naive Bayes, Support Vector Machines ou Régression Linéaire sont considérés comme des algorithmes classiques d'apprentissage automatique. Ils fonctionnent bien lorsqu'ils apprennent avec de petits jeux de données. Mais ils nécessitent souvent des lecteurs complexes. La tâche accomplie par les lecteurs est également appelée 'feature engineering'. Cela signifie qu'un être humain doit consacrer du temps à une analyse exploratoire approfondie du jeu de données.

Leurs caractéristiques peuvent être la fréquence des mots ou des lettres, mais aussi des éléments syntaxiques comme les noms, les adjectifs ou les verbes. Les caractéristiques les plus importantes pour la tâche à résoudre doivent être soigneusement sélectionnées et transmises à l'algorithme classique d'apprentissage automatique. Ce processus diffère de celui des réseaux de neurones. Lors de l'utilisation d'un réseau de neurones, il n'est pas nécessaire de recourir au 'feature engineering'. Les humains peuvent transmettre les données directement au réseau et obtiennent généralement de bonnes performances dès le départ. Cela permet d'économiser beaucoup de temps et de ressources.

L'inconvénient de la collaboration avec les réseaux de neurones est que vous avez besoin de beaucoup plus de données pour entraîner votre modèle de prédiction. Pensez à au moins 1 Go de fichiers texte. Pour vous donner une référence, 1 A4, soit un fichier texte de 5000 caractères, ne pèse que 5 Ko. Il vous faudrait donc 8.589.934 pages. Traiter plus de données sous-entend d'avoir accès à ces données et surtout, d'avoir beaucoup plus de puissance de traitement.

--- Les N-grammes de caractères pour la reconnaissance d'un auteur ---

Imaginez... vous travaillez pour une entreprise depuis plus de dix ans. Vous avez écrit des tonnes de courriels, d'articles, de notes internes et de rapports sur des sujets et dans des genres très différents. Tous vos écrits, ainsi que ceux de vos collègues, sont sauvegardés en toute sécurité sur les serveurs de l'entreprise.

Un jour, vous tombez amoureux d'une collègue. Après un certain temps, vous réalisez que cette personne est non seulement folle et hystérique mais qu'elle dépend beaucoup de vous. Le jour où vous décidez de rompre, votre ex élabore un plan pour vous tuer. Elle réussit. Pas de chance. Une lettre de suicide signée de votre nom est retrouvée à côté de votre cadavre. Celle-ci raconte que vous avez décidé de mettre fin à votre vie à cause de problèmes émotionnels. Vos meilleurs amis ne croient pas au suicide. Ils décident de porter l'affaire devant les tribunaux. Et là, à partir des textes que vous et d'autres avez produits, un modèle d'apprentissage automatique révèle que la

lettre de suicide a été écrite par quelqu'un d'autre.

Comment une machine analyse-t-elle les textes pour vous identifier ? La caractéristique la plus robuste pour la reconnaissance de l'auteur est fournie par la technique des N-grammes de caractères. Elle est utilisée dans des cas qui présentent une grande variété dans les thématiques et les genres d'écriture. Lors de l'utilisation des N-grammes de caractères, les textes sont considérés comme des séquences de caractères. Considérons le trigramme des caractères. Toutes les séquences de trois caractères qui se chevauchent sont isolées. Par exemple, le trigramme de caractères de 'suicide', serait, 'sui', 'uic', 'ici', 'cid' et 'ide'. Les N-grammes de caractères sont très simples, ils sont indépendants du langage et tolérants au bruit. De plus, les fautes d'orthographe ne compromettent pas la technique.

Les motifs trouvés avec les N-grammes de caractères se concentrent sur les choix stylistiques qui sont faits inconsciemment par l'auteur. Les modèles restent stables sur toute la longueur du texte, ce qui est important pour reconnaître l'auteur. D'autres types d'expériences pourraient inclure la longueur des mots ou des phrases, la richesse du vocabulaire, la fréquence des mots de fonction et même les mesures syntaxiques ou sémantiques.

Cela signifie non seulement que votre empreinte physique est unique, mais qu'il en va de même de la façon dont vous composez vos pensées !

La même technique n-gramme a découvert que 'The Cuckoo's Calling', un roman de Robert Galbraith, a en fait été écrit par... J.K. Rowling !

Références :

- Essai: On the Robustness of Authorship Attribution Based on Character N-gram Features, Efsthios Stamatatos, in Journal of Law & Policy, Volume 21, Issue 2, 2013.
- Article: <https://www.scientificamerican.com/article/how-a-computer-program-helped-show-jk-rowling-write-a-cuckoos-calling/>

--- Histoire des N-grammes ---

L'algorithme des N-grammes peut être retracé jusqu'aux travaux de Claude Shannon en théorie de l'information. Dans l'article 'A mathematical theory of communication', publié en 1948, Claude Shannon réalise la première instance d'un modèle de langage naturel à base des N-grammes. Il a posé la question suivante : étant donné la séquence des lettres, quelle est la probabilité de la prochaine lettre ? Si vous lisez l'extrait suivant, pouvez-vous nous dire par qui il a été écrit ? Shakespeare ou un robot N-grammes ?

SEBASTIEN : Dois-je rester debout jusqu'à la rupture.

BIRON : Cache ta tête.

VENTIDIUS : Il se rendit à Athènes, où, par le voeu. que j'ai fait pour m'occuper de toi.

FALSTAFF : Mon bon fripouille.

Vous aviez peut-être deviné, en considérant le sujet de ce récit, qu'un algorithme N-grammes a généré ce texte. Le modèle est entraîné sur l'oeuvre complète de Shakespeare. Alors que les algorithmes plus récents, tels que les réseaux de neurones récurrents de CharRNN, deviennent célèbres pour leurs performances, les N-grammes exécutent encore beaucoup de tâches NLP. Elles sont utilisées dans la traduction automatique, la reconnaissance vocale, la correction orthographique, la détection d'entités, l'extraction d'informations, etc.

Référence : <http://www.math.harvard.edu/~ctm/home/text/others/shannon/entropy/entropy.pdf>

--- Dieu dans Google Books ---

En 2006, Google crée un jeu de données de N-grammes à partir de sa collection de livres numérisés pour le mettre en ligne. Récemment, ils ont également réalisé une visionneuse de N-grammes.

Cela a permis de nombreuses recherches sociolinguistiques. Par exemple, en octobre 2018, le New York Times Magazine a publié un article d'opinion intitulé 'It's Getting Harder to Talk About God'.

L'auteur, Jonathan Merritt, avait analysé la mention du mot 'Dieu' dans le jeu de données de Google à l'aide du visualiseur de N-grammes. Il a conclu qu'il y a eu un déclin dans l'usage du mot depuis le 20ème siècle. Le corpus de Google contient des textes du 16e jusqu'au 21e siècle. Cependant l'auteur a manqué d'observer la popularité croissante des revues scientifiques vers le début du 20ème siècle. Ce nouveau genre, dans lequel le mot Dieu n'apparaît pas, a fait basculer le jeu des données. Si la littérature scientifique était retirée du corpus, la fréquence du mot 'Dieu' s'écoulerait toujours comme l'ondulation douce d'une vague lointaine.

Référence : <https://www.nytimes.com/2018/10/13/opinion/sunday/talk-god-spirituality-christian.html>

--- Les traits grammaticaux extraits de Twitter influencent le marché boursier ---

Les frontières entre les disciplines académiques s'estompent. La recherche économique mêlée à la psychologie, aux sciences sociales, aux concepts cognitifs et émotionnels créent un nouveau sous-domaine économique, appelé 'l'économie comportementale'.

Cela signifie que les chercheurs commencent à ex-

pliquer un mouvement boursier basé sur d'autres facteurs que les facteurs purement économiques. La Bourse et 'l'opinion publique' s'influencent mutuellement. De nombreuses recherches sont effectuées sur la façon d'utiliser 'l'opinion publique' pour prédire les tendances dans le cours des actions.

'L'opinion publique' est évaluée à partir de grandes quantités de données publiques, comme les tweets, les blogs ou la presse en ligne. Des recherches montrent que l'évolution des cours boursiers peut, dans une certaine mesure, être prédit en examinant 'l'opinion publique' à travers l'analyse des données automatique. On trouve de nombreux articles scientifiques en ligne, qui analysent la presse sur le 'sentiment' qui y est exprimé. Un article peut être annoté comme plus ou moins positif ou négatif. Les articles de presse annotés sont ensuite utilisés pour entraîner un modèle d'apprentissage automatique, qui permet de prédire les tendances boursières, en les marquant comme 'à la baisse' ou 'à la hausse'. Quand une entreprise fait mauvaise presse, les traders vendent. Au contraire, si les nouvelles sont bonnes, ils achètent.

Un article de Haikuan Liu de l'Université Nationale Australienne affirme que le temps des verbes utilisés dans les tweets peut être un indicateur de la fréquence des transactions financières. Son idée s'inspire du fait que la conjugaison des verbes est utilisée en psychologie pour détecter les premiers stades de la dépression humaine.

Référence : Grammatical Feature Extraction and Analysis of Tweet Text: An Application towards Predicting Stock Trends, The Australian National University (ANU)

--- Sac de mots ---

Dans le traitement du langage naturel, le 'sac de mots' est considéré comme un modèle simple. Il dépouille un texte de son contexte et le décompose dans sa collection de mots uniques. Ensuite, ces mots sont comptés. Dans les phrases précédentes, par exemple, le mot 'mots' est mentionné trois fois, mais ce n'est pas nécessairement un indicateur de l'objet du texte.

La première apparition de l'expression 'sac de mots' semble remonter à 1954. Zellig Harris a publié un article dans le contexte des études linguistiques, intitulé 'Distributional Structure'. Dans la partie intitulée 'Le sens en fonction de la distribution', il dit que 'le langage n'est pas seulement un sac de mots, mais aussi un outil aux propriétés particulières qui ont été façonnées au cours de son utilisation. Le travail du linguiste est précisément de découvrir ces propriétés, que ce soit pour l'analyse descriptive ou pour la synthèse du système quasi-linguistique.'

c us 'l8 t n | d cri i s ++++++ o ++++++ 3ini sst5 dl e or%tu ed5 u_
u1u t r éa n 86 Mi V |l|e|a|r|n|e|r|s| |l|e|a|r|n| / nV 8 r _ u s 9e 5mn ieo
a i8- | 3 se7 |s dtr ++++++ iT ++++++ 9l l o é e e r | _ _ ve- e
f ym it les3 63 d 9s5ue8 | s / , ,4 7 asr 9 d % 6sd il43 23 G |r o9c n % Dm u
i db o /mê eev oeut r +5o e s v9c + 7r ia 3 lii| p f H1 oo y L n7 eu
s + -0 i _ean, rsm 0 ++++++ s ++++++ ++++++ i qua u utn i + + l e
i_ t l 3 b% 7 éo s fn |l|e|a|r|n|e|r|s| r. |a|r|e| |p|a|t|t|e|r|n| rld 3s o l to N ,A 1 s
e s v l c ud ri50 ++++++ e ++++++ ++++++ %9e-am |e e l
rea r- 36 1- ru . 4 --ê e `t i ++++++ t c e a srn 9p 8 3 fm
eg u% 1k 6d gr- e s2 |f|i|n|d|e|r|s| , I u_o t s it n 8d d e
5,e - 9 et 2 A o 2e - e ae r l e ++++++ % n u +pa+ un8ne r
- r a ,n _ng |e r fhs ++++++ 2 ++++++ ++++++ e j h lqi n oae c7 r
6 t 9 % r e 6n ei enn |l|e|a|r|n|e|r|s| /i |a|r|e| |c|r|a|w|l|i|n|g| e 68 t 2c+ y t
5 u , 4 ue t e ++++++ a ++++++ ++++++ 6 n c5 b u45i u t| n
_ 7 i e l a 1 A78 it t ++++++ ++++++ 2 iid d ,a éc
0 n %t/ / h 4 r i6 7 sna ps |t|h|r|o|u|g|h| |d|a|t|a| e V lauae9n 5 9 l1
e ê Iho e t c t3 e 1 7 m ndm . 2 ++++++ ++++++ i ea s n p n c2 is
|\ 5 6 / ac r4 o ++++++ -ad6 ++++++ ++++++ e3 cv 8
s u% - t e 1 3o |l|e|a|r|n|e|r|s| e m |g|e|n|e|r|a|t|e| |s|o|m|e| |k|i|i|n|d| t L
s 4rr s t - us l ++++++ / | ++++++ ++++++ ae r
g '- ê u u 2 -4 48 9s ++++++ ++++++ m s +i 39 +7p o - di
o t2 4 e l a os b i o |o|f| |s|p|e|c|i|f|i|c| f -p ee 24 7 e
r -e o lr 3i dh t t Co ++++++ ++++++ ee3 r 4 r | ind
id ei t e 7m 6 g t p + é e e ++++++ ++++++ s _a ni i l t w ' o
r c 7 s_ l ea t ' - |'g|r|a|m|m|a|r|' | o 0 t 1 u n uf -
teu| c_n e o5 e e % d+ re 5 t ++++++ ++++++ n s o etps r ma iu a
r u s 5 ++++++ ++++++ ++++++ 2ain d%
p 7,s0 .e e 1r gt _ |c|l|a|s|s|i|f|i|e|r|s| ir |g|e|n|e|r|a|t|e|, |e|v|a|l|l|u|a|t|e| ê sa, g
te (i m e + ++++++ a ++++++ ++++++ es a _ s
e dr p ' o U + p' ++++++ ++++++ 1 o n8 - r 8
r6 o é s re eA a s r' n c t 5 |a|n|d| |r|e|a|d|j|u|s|t| o u r + n2o t o
t % 4 iinu g 4 7s p u u- e 7 ++++++ ++++++ ' u tn + - e
l/ v 2 ni 2 a ++++++ u ++++++ ++++++ 35 8 _ p
_ %o dh 9 ee i s |l|e|a|r|n|e|r|s| 3a |u|n|d|e|r|s|t|a|n|d| |a|n|d| s ee c i %
m gra st ++++++ h ++++++ ++++++ d+ g a n
p seh . t r - d si r- i i ++++++ ++++++ ++++++ n 6 9 r
s e i i 2 t t ae i |r|e|v|e|a|l| |p|a|t|t|e|r|n|s| me ' eg
e i - tB i \m t ++++++ ++++++ 1 t e 45
5 s u iee i i ,l 7t ++++++ ++++++ 43 l u .a d r
n 59n e l w 8a |l|e|a|r|n|e|r|s| 8 |d|o|n|'t| |a|l|w|a|y|s| l i i d r
l è p m r 0 i ++++++ % ++++++ ++++++ 4 i ee t 5
m -8 p t a 0 en v6 ++++++ ++++++ p é t
di ' i 2 '4 uf e c l t |d|i|s|t|u|i|n|g|u|i|s|h| |w|e|l|l|l| m l e+
n5 de .e 2 r _ ++++++ ++++++ z 1 - d
_ , 2 l csi d |w|h|i|c|h| |p|a|t|t|e|r|n|s| r i a.
v a l 9 i s 9 sk ++++++ ++++++ 3 a
3 e s r A e s 9 r al ++++++ ++++++ ++++++ hsé
r u ue I C _ 4 m ++++++ ++++++ ++++++ +
e a n i v 6s e /a 5 7C r ' r64 -
a l e ei a _ 6 e nu l n ' r64 -
roe l e e 2 6e n a 9
o g . e - /q 2 m 7 .1 1 |n . q 8 9p s7
h o u fe 6 r n a ê n 5 4
- nn n i u . r dt | r t 4V
9 7 n l s e - i 9n 9
4 n m 8 r n a_ 4
1 s b , ui + % e \ e e s
s | _ 0 i il s a
6e a + e d 9 p t a t
e e ' r i 9 t 9 _ n i s
- n r s7 s e 9 m + s
è e Q 8e . t s5 eo t c b ' 1
% e 0 s _ b n n 4
% v e ' it 2 r u r
o r pu p r

--- Naive Bayes & Viagra ---

L'algorithme Naive Bayes est un Apprenant célèbre qui réussit bien avec peu de données. Nous l'appliquons tout le temps. Christian & Griffiths affirment dans leur livre, 'Algorithms to Live by', que 'nos jours sont remplis de petites données'. Imaginez par exemple que vous vous trouviez à un arrêt de bus dans une ville étrangère. L'autre personne qui se tient là attend depuis 7 minutes. Qu'est-ce que vous faites ? Décidez-vous d'attendre ? Et si oui, pour combien de temps ? Quand allez-vous envisager d'autres options ? Un autre exemple. Imaginez qu'un ami demande conseil sur une relation. Il est avec son nouveau partenaire depuis un mois. Doit-il l'inviter à l'accompagner à un mariage de famille ?

Les croyances préexistantes sont cruciales pour que Naive Bayes fonctionne. L'idée est de calculer les probabilités sur base de ces connaissances préalables et d'une situation spécifique.

Le théorème a été formulé dans les années 1740 par le révérend et mathématicien amateur Thomas Bayes. Il a consacré sa vie à résoudre la question de savoir comment gagner à la loterie. Mais la règle de Bayes a été rendue célèbre dans sa forme actuelle par le mathématicien Pierre-Simon Laplace en

temps après la mort de La Place, la théorie tombe dans l'oubli jusqu'à ce qu'elle soit à nouveau déterrée pendant la Seconde Guerre mondiale dans le but de briser le code Enigma.

La plupart des personnes sont aujourd'hui entrées en contact avec Naive Bayes par le biais de leurs dossiers de courrier indésirable. Naive Bayes est un algorithme largement utilisé pour la détection du spam. C'est une coïncidence que le Viagra, médicalement contre la dysfonction érectile, a été approuvé par la FDA (US Food & Drug Administration) en 1997, au moment où environ 10 millions d'utilisateurs dans le monde avaient des comptes de messagerie Web gratuits. Les sociétés de vente avaient l'intelligence d'utiliser la publicité massive par e-mail : c'était un média intime, à l'époque réservé à la communication privée. En 2001, le premier programme SpamAssasin s'appuyant sur Naive Bayes a été téléchargé sur SourceForge, réduisant ainsi le marketing 'guerilla par courriel'.

Référence : Machine Learners, by Adrian MacKenzie, The MIT Press, Cambridge, US, November 2017.

--- Naive Bayes & Enigma ---

Cette histoire de Naive Bayes fait partie du livre 'The theory that would not die', écrit par Sharon Bertsch McGrayne. Elle décrit entre autres comment Naive Bayes est vite oubliée après la mort de

Pierre-Simon Laplace, son inventeur. Le mathématicien aurait échoué à créditer les travaux des autres. Par conséquent, il a souffert d'accusations largement diffusées contre sa réputation. Ce n'est que 150 ans plus tard que l'accusation s'est avérée fausse.

Avançons en 1939, alors que le règne de Bayes demeure pratiquement tabou, mort et enterré dans le domaine de la statistique. Lorsque la France est occupée en 1940 par l'Allemagne, qui contrôle les usines et les fermes européennes, la plus grande inquiétude de Winston Churchill est le péril U-boot. Les opérations de sous-marin étaient étroitement contrôlées par le quartier général allemand en France. Chaque sous-marin partait en mer sans ordres, et les recevait sous forme de messages radio codés après avoir atteint l'Atlantique. Les messages étaient cryptés par des machines à brouiller les mots, appelées Enigma machines. Enigma ressemblait à une machine à écrire compliquée. Elle est inventée par la société allemande Scherbius & Ritter après la première guerre mondiale, lorsque le besoin de machines d'encodage de messages est devenu douloureusement évident.

Curieusement, et heureusement pour Naive Bayes et le monde, à l'époque le gouvernement britannique et les systèmes d'éducation considéraient les mathématiques appliquées et les statistiques sans aucun rapport avec la résolution pratique des problèmes. Les données statistiques ont été jugées gênantes en raison de leur caractère détaillé. Ainsi, les données du temps de guerre étaient souvent analysées non pas par des statisticiens, mais par des biologistes, des physiciens et des mathématiciens théoriques. Aucun d'entre eux ne savait qu'en ce qui concerne les statistiques sophistiquées, la règle de Bayes était considérée non-scientifique.

C'est le désormais célèbre Alan Turing, mathématicien, informaticien, logicien, cryptanalyste, philosophe et biologiste théorique, qui a utilisé le système de probabilités des règles de Bayes pour concevoir la 'bombe'. Il s'agissait d'une machine électromécanique à grande vitesse pour tester tous les arrangements possibles qu'une machine Enigma produirait. Afin de déchiffrer les codes navals des U-boot, Turing simplifie le système de la 'bombe' en utilisant des méthodes bayésiennes. La 'bombe' a transformé le quartier général du Royaume-Uni en une usine de décryptage. L'histoire est bien illustrée dans 'The Imitation Game', un film de Morten Tyldum, sorti en 2014.

--- Une histoire sur les petits pois ---

En statistique, la régression linéaire est une méthode d'apprentissage supervisé. Après l'entraînement avec des données annotées, le modèle tente de prédire les valeurs de nouvelles données inconnues. La régression linéaire permet de résumer et

d'étudier les relations entre deux éléments, afin de voir s'il existe une corrélation entre eux.

S'il y a une corrélation positive, la connaissance d'un élément aide à prédire l'autre. Par exemple, étant donné la critique d'un film, nous pouvons prédire le nombre moyen d'étoiles qui lui sont attribuées, plutôt que de simplement dire si la critique est positive ou négative.

Parfois, les figures que nous rencontrons en grattant sous la surface ne sont pas à notre goût. L'idée de régression vient de Sir Francis Galton, un scientifique influent du 19e siècle. Il a passé sa vie à étudier le problème de l'hérédité - pour comprendre à quel point les caractéristiques d'une génération d'êtres vivants se manifestent dans la génération suivante. Il a établi le domaine de l'eugénisme et l'a défini comme 'l'étude des organismes sous contrôle social qui peuvent améliorer ou altérer les qualités raciales des générations futures, que ce soit physiquement ou mentalement'. Par conséquent, son nom a marqué l'histoire et l'héritage du racisme scientifique.

Galton a d'abord abordé le problème de l'hérédité en examinant les caractéristiques du petit pois doux. Il a choisi le petit pois parce que l'espèce peut s'auto-fertiliser. Les plantes femelles héritent des variations génétiques des plantes mères sans la contribution d'un deuxième parent. Cette caractéristique élimine la nécessité de traiter avec des sources multiples.

En 1875, Galton a distribué des paquets de graines de petits pois à sept amis. Chaque ami recevait des graines de poids uniforme, mais il y avait des variations importantes d'un paquet à l'autre. Les amis de Galton ont récolté les graines des nouvelles générations de plantes et les lui ont rendues. Il a ensuite tracé le poids des graines femelles contre le poids des graines mères. Il a découvert que le poids médian des graines femelles d'une taille particulière de la semence mère décrivait approximativement une ligne droite avec une pente positive inférieure à 1,0. Les premières idées de Galton sur la régression sont nées de ce diagramme bidimensionnel qui compare la taille des petits pois femelles à celle des petits pois mères. Il a utilisé cette représentation de ses données pour illustrer les fondements de ce que les statisticiens appellent encore aujourd'hui la régression. Pour Galton, c'était aussi une façon de décrire les avantages de l'eugénisme.

La recherche de Galton a été appréciée par de nombreux intellectuels de son temps. En 1869, dans 'Hereditary Genius', Galton affirme que le génie est principalement une question d'ascendance. Il croyait qu'il y avait une explication biologique à l'inégalité sociale entre les races. Galton a même persuadé son demi-cousin Charles Darwin de ses idées. Après avoir lu l'article de Galton, Darwin a déclaré : 'Vous avez converti un adversaire, car

j'ai toujours soutenu qu'à l'exception des imbéciles, les hommes ne différaient pas beaucoup sur le plan intellectuel, seulement sur le plan du zèle et du labeur'. Heureusement, l'étude moderne de l'hérédité a réussi à éliminer le mythe de la différence génétique fondée sur la race.

La raison pour laquelle nous l'évoquons dans cette série, c'est qu'il a été parmi les premiers scientifiques à utiliser des méthodes statistiques dans ses recherches. Sa principale contribution dans ce domaine a été l'analyse de régression linéaire, qui a fondé les bases d'une grande partie de la statistique moderne. Alors que nous nous engageons dans le domaine de l'apprentissage automatique, Algolite essaie de ne pas oublier que les systèmes d'ordre ont du pouvoir, et que ce pouvoir n'a pas toujours été exercé au bénéfice de tout le monde. L'apprentissage automatique a hérité de nombreux aspects de la recherche statistique, certains plus agréables que d'autres. Nous devons nous méfier, car ces visions du monde s'infiltrèrent dans les modèles algorithmiques qui créent des ordres aujourd'hui.

Références :

<http://galton.org/letters/darwin/correspondence.htm>

<https://www.tandfonline.com/doi/full/10.1080/10691898.2001.11910537>

<http://www.paramoulipist.be/?p=1693>

--- Perceptron ---

Nous nous trouvons dans une décennie où les réseaux de neurones suscitent beaucoup d'attention. Cela n'a pas toujours été le cas. L'étude des réseaux de neurones remonte aux années 1940, lorsque la première métaphore des neurones est apparue. Le neurone n'est pas la seule référence biologique dans le domaine de l'apprentissage automatique - pensez au mot corpus ou formation. Le neurone artificiel a été construit en relation étroite avec son homologue biologique.

Le psychologue Frank Rosenblatt s'est inspiré des travaux de son collègue Donald Hebb sur le rôle des neurones dans l'apprentissage humain. Hebb a déclaré que 'les cellules qui communiquent, se mettent ensemble.' Sa théorie est maintenant à la base de l'apprentissage associatif humain, mais aussi de l'apprentissage en réseau de neurones non supervisé. Il a poussé Rosenblatt à développer l'idée du neurone artificiel. En 1962, il crée le Perceptron. Le Perceptron est un modèle qui apprend par la pondération des entrées.

Il a été mis de côté par les chercheurs, parce qu'il ne peut gérer que la classification binaire. Cela signifie que les données doivent être sépa-

rables linéairement, comme par exemple hommes et femmes, noir et blanc. Il est clair que ce type de données est très rare dans le monde réel. Lorsque le soi-disant premier hiver de l'Intelligence Artificielle (IA) est arrivé en 1974-1980 et que le financement consacré à cette recherche a diminué, le Perceptron a également été négligé. Pendant 10 ans, il est resté inactif. Lorsque le printemps s'installe à la fin des années 1980, de nouvelles générations de chercheurs le reprennent et l'utilisent pour construire des réseaux de neurones. Ceux-ci contiennent de multiples couches de Perceptrons. C'est ainsi que les réseaux de neurones voient la lumière. On pourrait dire que cette saison d'apprentissage automatique est particulièrement chaude, mais il faut un autre hiver pour connaître un été.

--- BERT ---

Certains articles en ligne disent que l'année 2018 a marqué un tournant dans le domaine du traitement du langage naturel. Une série de modèles de 'deep learning' ont permis d'obtenir des résultats excellents pour des tâches comme les réponses aux questions ou la classification des sentiments. L'algorithme BERT de Google est entré dans les concours d'apprentissage automatique de l'année dernière comme un 'modèle gagnant'. Il témoigne d'une performance supérieure sur une grande variété de tâches.

BERT est pré-entraîné; ses poids sont appris à l'avance grâce à deux tâches non supervisées. Cela signifie que BERT n'a pas besoin d'être entraîné à partir de zéro pour chaque nouvelle tâche. Vous n'avez qu'à affiner ses poids.

Cela signifie également qu'un programmeur souhaitant utiliser BERT ne sait plus sur quels paramètres BERT est réglé, ni à base de quelles données il a appris ses performances.

BERT signifie 'Bidirectional Encoder Representations from Transformers'. Cela signifie que BERT permet un entraînement bidirectionnel. Le modèle apprend le contexte d'un mot à partir de son environnement, à gauche et à droite d'un mot. En tant que tel, il peut faire la différence entre 'Je suis pile à l'heure' et 'Je l'ai mis sur la pile'.

Quelques faits :

- BERT_large, avec 345 millions de paramètres, est le plus grand modèle du genre. Il est manifestement supérieur à BERT_base, qui utilise la même architecture avec 'seulement' 110 millions de paramètres, pour les tâches à petite échelle.

- Pour exécuter BERT, vous devez utiliser les TPU. Ce sont les processeurs (CPU) de Google spécialement conçus pour TensorFlow, la plateforme de

'deep learning'. Les tarifs de location de TPU vont de de 8\$/h à 394\$/h. Si vous êtes comme nous, et vous ne voulez pas travailler avec des solutions prêtes à l'emploi, et vous souhaitez ouvrir la boîte noire, BERT exige de faire des économies pour pouvoir l'utiliser.

Références :

<https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html>

<https://towardsdatascience.com/deconstructing-bert-distilling-6-patterns-from-100-million-parameters-b49113672f77>

0 12 3 4 5 67 8 9 0
 12 3 4 5 67 8 9 0 12
 3 4 5 67 8 9 0 1 2 3
 4 56 7 8 9 01 2 3
 4 56 7 8 9 01 2 3 4
 5 6 7 8 9 0 1 2 3 4 5 6
 7 8 9 0 1 2 3 4 5 6 7 8 9
 7 89 0 1 2 34 5 6 7 89
 89 0 1 2 3 4 5 6 7 8 9
 0 1 23 4 5 6 78 9 0
 1 23 4 5 6 78 9 0
 1 2 3 4 5 6 7 8 9 0 12
 3 4 5 67 8 9 0 12 3
 4 5 6 7 8 9 0 1 2 3
 4 56 7 8 9 01 2 3 4
 5 6 7 8 9 0 1 2 3 4 5 6
 7 8 9 0 1 2 3 4 5 6 7 8 9
 7 8 9 0 1 2 3 4 5 6 7 89
 89 0 1 2 34 5 6 7 89
 0 1 2 3 4 5 6 7 8 9 0
 1 23 4 5 6 78 9 0
 1 2 3 4 5 6 7 8 9 0 12 3
 4 5 67 8 9 0 12 3
 4 5 67 8 9 0 1 2 3
 4 5 6 7 8 9 01 2 3 4
 56 7 8 9 01 2 3 4 5
 6 7 8 9 0 1 2 3 4 5 6
 7 8 9 0 1 2 3 4 5 6 7 89
 7 8 9 0 1 2 3 4 5 6 7 89
 0 1 2 34 5 6 7 89 0
 1 2 34 5 6 7 8 9 0
 1 23 4 5 6 78 9 0
 1 23 4 5 6 7 8 9 0 12 3
 2 3 4 5 6 7 8 9 0 12 3
 4 5 67 8 9 0 12 3

GLOSSAIRE

Vous trouverez ci-dessous un glossaire non-exhaustif reprenant des termes fréquemment utilisés dans l'exposition. Il est conçu comme une aide pour les visiteurs connaissant peu le vocabulaire lié au domaine du traitement des langues naturelles (NLP), Algolit ou le Mundaneum.

* ALGOLIT

Un groupe bruxellois spécialisé dans la recherche artistique sur les algorithmes et la littérature. Chaque mois, le groupe se réunit pour expérimenter avec du code et des textes publiés sous licences libres.
<http://www.algolit.net>

* ALGOLITTÉRAIRE

Terme inventé par Algolit pour des œuvres qui explorent le point de vue du conteur algorithmique. Quelles nouvelles formes de narration rendons-nous possibles en dialoguant avec les algorithmes ?

* ALGORITHME

Un ensemble d'instructions dans un langage de programmation spécifique, qui permettent de produire un résultat (output) à partir de données (inputs).

* ANNOTATION

Le processus d'annota-

tion est une étape cruciale de l'apprentissage automatique supervisé durant laquelle l'algorithme reçoit des exemples de ce qu'il doit apprendre. Un filtre anti-spam sera alimenté d'exemples de messages spams et de messages réels. Ces exemples consistent en un message, l'entrée, accompagné d'une étiquette spam ou non spam. L'annotation d'un jeu de données est un travail exécuté par des humains, qui choisissent une étiquette pour chaque élément du jeu de données. Pour assurer la qualité des étiquettes, plusieurs annotateurs doivent voir le même élément, la même entrée, et donner la même étiquette avant qu'un exemple ne soit inclus dans les données d'entraînement.

* APPRENTISSAGE

AUTOMATIQUE OU MACHINE LEARNING
Modèles algorithmiques basés sur la statistique, principalement utilisés pour analyser et prédire des situations à partir de cas existants. Dans cette exposition, nous nous concentrons sur les modèles d'apprentissage automatique pour le traitement de texte ou le traitement du langage naturel (voir NLP). Ces modèles ont appris à effectuer une tâche spécifique sur la base de textes existants. Ils sont utilisés par les moteurs de recherche, les traductions automatiques, et permettent de générer des résumés et de repérer les tendances sur les réseaux sociaux et des fils d'actualité. Ils influencent ce que l'on

voit en tant qu'utilisateur, mais ont aussi leur mot à dire dans les fluctuations du cours des bourses mondiales ou dans la détection de la cybercriminalité et du vandalisme.

* APPRENTISSAGE AUTOMATIQUE CLASSIQUE

Naïve Bayes, Support Vector Machines ou Régression Linéaire sont considérés comme des algorithmes classiques d'apprentissage automatique. Ils fonctionnent bien lorsqu'ils apprennent avec de petits jeux de données. Mais ils nécessitent souvent des lecteurs complexes. La tâche accomplie par les lecteurs est également appelée 'feature engineering' (voir ci-dessous). Cela signifie qu'un être humain doit consacrer du temps à une analyse exploratoire approfondie du jeu de données.

* BAG OF WORDS

Le modèle du sac de mots est une représentation simplifiée du texte utilisé dans le traitement du langage naturel. Dans ce modèle, un texte est représenté sous la forme d'une collection de mots uniques, sans tenir compte de la grammaire, de la ponctuation ni même de leur ordre dans le texte. Ce modèle transforme un texte en une liste de mots associés à leur fréquence littéralement un sac de mots. Le sac de mots est souvent utilisé comme référence, c'est sur cette base qu'on évaluera la performance d'un nouveau modèle.

* CHAÎNE DE MARKOV

Algorithme qui scanne un texte à la recherche de la probabilité de tran-

sition d'occurrences de lettres ou de mots, ce qui donne des tables de probabilité de transition qui peuvent être calculées sans aucune compréhension sémantique ou grammaticale du langage naturel. Cet algorithme peut être utilisé pour analyser des textes, mais aussi pour les recombiner. Il est largement utilisé pour la génération de spam.

* CONSTANT

Constant est une association sans but lucratif d'artistes autogérés, basée à Bruxelles depuis 1997 et active dans les domaines de l'art, des médias et de la technologie. Algolit est né en 2012 comme un projet de Constant.
<http://constantvzw.org>

* DATA WORKERS

Intelligences artificielles développées pour servir, divertir, enregistrer et connaître les humains. Le travail de ces entités machiniques est généralement dissimulé derrière des interfaces et des brevets. Dans l'exposition, les conteurs algorithmiques quittent leur monde souterrain invisible pour devenir nos interlocuteurs.

* DONNÉES D'ENTRAÎNEMENT

Les algorithmes d'apprentissage automatique ont besoin d'être guidés. Pour séparer une chose d'une autre, faire des distinctions, ils ont besoin de motifs. Ils les trouvent dans les textes qui leur sont donnés, les données d'entraînement. L'être humain doit choisir avec soin un matériel d'entraînement adapté à la tâche de la machine. Il n'est pas logique d'en-

traîner une machine avec des romans du 19ème siècle si sa mission est d'analyser des Tweets.

* DUMP

Terme anglais signifiant 'dépôt, décharge, déverser massivement'. En informatique, le terme dump désigne généralement une copie brute d'une base de données; par exemple pour effectuer une sauvegarde de données ou pour les utiliser ailleurs. Les dumps sont souvent publiés par des projets de logiciels libres et de contenu libre, tels que Wikipédia, pour permettre la réutilisation ou la dérivation(fork) de la base de données.

* FEATURE ENGINEERING

Processus utilisant la connaissance du domaine des données pour créer les caractéristiques qui font fonctionner les algorithmes d'apprentissage machine. En d'autres termes, un être humain doit consacrer du temps à une analyse exploratoire approfondie du jeu de données, afin d'en définir les principales caractéristiques. Ces caractéristiques peuvent être la fréquence des mots ou des lettres, mais aussi des éléments syntaxiques comme les noms, les adjectifs ou les verbes. Les caractéristiques les plus importantes pour la tâche à résoudre doivent être soigneusement sélectionnées pour être transmises à un algorithme classique d'apprentissage automatique.

* FLOSS OU LOGICIELS LIBRES ET OPEN SOURCE

Un logiciel libre est un logiciel dont l'utilisation, l'étude, la modification et la duplica-

tion par autrui en vue de sa diffusion sont permises, techniquement et légalement, ceci afin de garantir certaines libertés induites, dont le contrôle du programme par l'utilisateur et la possibilité de partage entre individus. Ces droits peuvent être simplement disponibles - cas du domaine public - ou bien établis par une licence, dite 'libre', basée sur le droit d'auteur. Les 'licences copyleft' garantissent le maintien de ces droits aux utilisateurs même pour les travaux dérivés. Les logiciels libres constituent une alternative à ceux qui ne le sont pas, qualifiés de 'propriétaires' ou de 'privatiseurs'. (Wikipedia)

* GIT

Un système logiciel permettant de suivre les changements dans le code source pendant le développement d'un logiciel. Il est conçu pour coordonner le travail des programmeurs, mais il peut être utilisé pour suivre les changements dans n'importe quel ensemble de fichiers. Avant d'initier un nouveau projet, les programmeurs créent un 'dépôt git' dans lequel ils publieront toutes les parties du code. Les dépôts git d'Algolit se trouvent ici <https://gitlab.constantvzw.org/algolit>.

* GUTENBERG.ORG

Le projet Gutenberg est une bibliothèque de versions électroniques libres de livres physiquement existants. Les textes fournis sont essentiellement du domaine public, soit parce qu'ils n'ont jamais été

sujets à des droits d'auteur soit parce que ces derniers sont expirés. Le projet fut lancé par Michael Hart en 1971 et nommé en hommage à l'imprimeur allemand du XVe siècle Johannes Gutenberg. (Wikipedia)

* HENRI LA FONTAINE

Henri La Fontaine (1854-1943) est un homme politique, féministe et pacifiste belge. Il reçoit le Prix Nobel de la paix en 1913 en raison de son engagement au sein du Bureau International de la Paix et de sa contribution à l'organisation du mouvement pacifiste. En 1895, ensemble avec Paul Otlet, il crée ensemble l'Institut international de bibliographie qui deviendra le Mundaneum. Au sein de cette institution, qui visait à rassembler l'ensemble des connaissances du monde, il contribue à mettre au point le système de Classification décimale universelle (CDU).

* IA OU INTELLIGENCES ARTIFICIELLES

L'intelligence artificielle (IA) est 'l'ensemble des théories et des techniques mises en œuvre en vue de réaliser des machines capables de simuler l'intelligence. Elle correspond donc à un ensemble de concepts et de technologies plus qu'à une discipline autonome constituée. D'autres, remarquant la définition peu précise de l'IA, notamment la CNIL, la définissent comme 'le grand mythe de notre temps'. (Wikipedia)

* KAGGLE

Plateforme en ligne où les utilisateurs trouvent et publient des ensembles de données,

explorent et construisent des modèles d'apprentissage automatique, collaborent avec d'autres et participent à des concours pour relever des défis. Environ un demi-million d'utilisateurs sont actifs sur Kaggle. Kaggle a été fondée par Goldbloom et Ben Hamner en 2010 et acquise par Google en mars 2017.

* LANGAGE NATUREL

Selon Wikipédia, 'Une langue dite « naturelle » est une langue qui s'est formée petit à petit, évoluant avec le temps, et qui fait partie du langage naturel. Son origine est bien souvent floue et peut être retracée plus ou moins clairement par la linguistique comparée. On oppose les langues naturelles - comme le français - aux langues construites comme le langage de programmation ou l'espéranto, formées intentionnellement par l'entremise de l'homme pour remplir un besoin précis.'

* LITTÉRATURE

Algolit comprend la notion de littérature comme beaucoup d'autres auteurs expérimentaux elle inclut toute la production linguistique, du dictionnaire à la Bible, de l'œuvre entière de Virginia Woolf à toutes les versions des Conditions d'utilisation publiées par Google depuis son existence. En ce sens, le code de programmation peut aussi être de la littérature.

* MECHANICAL TURK

Le Mechanical Turk d'Amazon est une plateforme en ligne à destination des humains conçue pour exécuter des tâches que

les algorithmes ne parviennent pas à faire. Il peut s'agir, par exemple, d'annoter des phrases comme étant positives ou négatives, de repérer des plaques d'immatriculation, de reconnaître des visages. Les annonces que l'on trouve sur cette plateforme sont souvent rémunérées moins d'un centime par tâche. Les tâches les plus complexes ou nécessitant le plus de connaissances peuvent être payées jusqu'à plusieurs centimes. De nombreux chercheurs universitaires utilisent le Mechanical Turk pour des tâches qui auraient été exécutées par des étudiants auparavant.

* MODÈLES D'APPRENTISSAGE AUTOMATIQUE SUPERVISÉ

Pour la création de modèles d'apprentissage automatique supervisés, les humains annotent les échantillons d'entraînement avant de les envoyer à la machine. Chaque texte est jugé par au moins 3 humains par exemple, s'il s'agit de spam ou non, s'il est positif ou négatif.

* MODÈLES D'APPRENTISSAGE AUTOMATIQUE NON-SUPERVISÉ

Les modèles d'apprentissage automatique non supervisés n'ont pas besoin de l'étape d'annotations des données par des humains. Par contre, ils nécessitent de grandes quantités de données pour s'entraîner.

* MUNDANEUM

À la fin du 19ème siècle, deux jeunes juristes belges, Paul Otlet (1868-1944), 'père de la documentation', et Henri La Fontaine

(1854-1943), homme d'État et prix Nobel de la paix, créent le Mundaneum. Le projet vise à rassembler toute la connaissance du monde et à la classer à l'aide du système de Classification décimale universelle (UDC) qu'ils inventent.

* NATURAL LANGUAGE PROCESSING (NLP)

Le traitement du langage naturel (NLP) est un terme collectif qui désigne le traitement informatique automatique des langues humaines. Cela comprend les algorithmes utilisant, comme données, du texte produit par l'homme et qui tentent de le reproduire.

* N-GRAMMES DE CARACTÈRES

une technique utilisée pour la reconnaissance de la paternité d'une oeuvre. Lors de l'utilisation des N-grammes de caractères, les textes sont considérés comme des séquences de caractères. Considérons le trigramme des caractères. Toutes les séquences de trois caractères qui se chevauchent sont isolées. Par exemple, le trigramme de caractères de 'suicide', serait, 'sui,' uic', uic', 'ici', 'cid', etc. Les motifs trouvés avec les N-grammes de caractères se concentrent sur les choix stylistiques qui sont faits inconsciemment par l'auteur. Ces modèles restent stables sur toute la longueur du texte.

* ORACLE

Les Oracles sont un type particulier de modèles algorithmiques souvent basés sur la statistique, qui servent à pré-

dire des situations particulières ou à profiler des habitudes d'utilisateurs. Elles sont largement utilisées dans les smartphones, les ordinateurs et les tablettes.

* OULIPO

Le collectif Oulipo, acronyme d'Ouvroir de Littérature Potentielle, est une grande source d'inspiration pour Algorithmic. Oulipo a été créé à Paris par les écrivains Raymond Queneau et François Le Lionnais. Ils ont ancré leur pratique dans l'avant-garde européenne du XXe siècle et dans la tradition expérimentale des années 60. Pour Oulipo, la création de règles devient la condition permettant de générer de nouveaux textes, ou ce qu'ils appellent la littérature potentielle. Plus tard, en 1981, ils ont également créé ALAMO - Atelier de Littérature Assistée par la Mathématique et les Ordinateurs.

* PAUL OTLET

Paul Otlet (1868 - 1944) était un auteur, entrepreneur, visionnaire, avocat et militant pour la paix belge ; il est l'une des nombreuses personnes qui ont été considérées comme le père des sciences de l'information, un domaine qu'il a appelé 'la documentation'. Otlet a créé la Classification décimale universelle, qui s'est répandue dans les bibliothèques. Avec Henri La Fontaine, il crée le Palais Mondial, qui devient le Mundaneum, pour abriter les collections et les activités de leurs différents organismes et instituts.

* PYTHON

Le principal langage de programmation utilisé dans le monde entier pour le traitement du langage, inventé en 1991 par le programmeur néerlandais Guido Van Rossum.

* RECONNAISSANCE OPTIQUE DE CARACTÈRES (ROC)

en anglais optical character recognition (OCR), ou océrisation, désigne les procédés informatiques permettant la traduction d'images de textes scannés en fichiers de texte manipulables.

* RÉSEAUX DE NEURONES

Systèmes informatiques inspirés des réseaux neuronaux biologiques trouvés dans le cerveau des animaux. Un réseau de neurone n'est pas un algorithme, mais plutôt un cadre dans lequel de nombreux algorithmes d'apprentissage machine différents travaillent ensemble et traitent des données complexes. De tels systèmes 'apprennent' à exécuter des tâches en observant des exemples, généralement sans être programmés à priori avec des règles spécifiques. Par exemple, un algorithme de reconnaissance de chat apprendra à identifier les images qui contiennent des chats en observant des images qui ont été étiquetées manuellement comme 'chat' ou 'pas chat'. Il utilisera ces exemples pour générer ce qu'il considère être un chat et pourra identifier les chats dans d'autres images. Il le fera sans aucune connaissance préalable sur les chats. Il générera automatiquement ses propres caractéristiques d'identi-

cation à partir du matériel d'apprentissage qui lui est donné.

* RULE-BASED MODELS

Les Oracles peuvent être créés à l'aide de différentes techniques. L'une d'entre elles consiste à définir manuellement les règles. Ces modèles sont appelés 'rule-based models' (modèles basés sur des règles), et se situent à l'opposé des modèles statistiques. Ils sont utiles pour des tâches spécifiques, comme par exemple, la détection de la mention d'une certaine molécule dans un article scientifique. Ils sont performants, même avec très peu de données d'entraînement.

* SENTIMENT ANALYSIS

Également appelé 'opinion mining' (sondage d'opinion). Une tâche fondamentale de l'analyse des sentiments consiste à classer un texte donné comme positif, négatif ou neutre. La classification avancée des sentiments 'au-delà de la polarité' examine, par exemple, les états émotionnels tels que 'en colère', 'triste' et 'heureux'. L'analyse du sentiment est largement appliquée aux actions des utilisateurs tels que les critiques et les réponses aux enquêtes, les commentaires et les messages sur les médias sociaux, et les documents de santé. Elle est intégrée dans des applications qui vont du marketing au service à la clientèle, des transactions boursières à la médecine clinique.

* TF-IDF (TERM FREQUENCY - INVERSE DOCUMENT FREQUENCY)

Une méthode de pondération utilisée dans la recherche de textes. Cette mesure statistique permet d'évaluer l'importance d'un terme contenu dans un document, relativement à une collection ou un corpus de textes. Le poids augmente proportionnellement au nombre d'occurrences du mot dans le document. Il varie également en fonction de la fréquence du mot dans le corpus. Le TF-IDF est notamment utilisé dans la classification des spams.

* 'WORD EMBEDDINGS'

Techniques de modélisation du langage qui, par de multiples opérations mathématiques, tracent des mots dans un espace vectoriel multidimensionnel. Lorsque les mots sont 'embedded' ou intégrés, ils se transforment de symboles distincts en objets mathématiques, qui peuvent être multipliés, divisés, ajoutés ou soustraits.

* WORDNET

Wordnet est une combinaison d'un dictionnaire et d'un thésaurus qui peut être lu par des machines. Selon Wikipédia, il a été créé dans le Cognitive Science Laboratory de l'Université de Princeton à partir de 1985.

