

◆ Thème 3.5. L'intelligence artificielle

I. — Traitement automatique de l'information

1) Les origines

Les premières machines à effectuer automatiquement une tâche donnée sont d'une part des « machines à calculer », ancêtre des calculatrices et, d'autre part, des métiers à tisser.

- 1642-1645 • Invention de la Pascaline, première « machine à calculer » par Pascal.
- 1671-1694 • Leibniz construit une machine pouvant réaliser des multiplications.
- 1725 • Bouchon invente un métier à tisser programmable à l'aide d'un ruban.
- 1728 • Falcon remplace le ruban par des cartes perforées.
- 1745-1755 • Vaucanson remplace les cartes perforées par un cylindre métallique.
- 1801 • Jacquard crée son métier à tisser reprenant les machines précédentes

Au début du XIXe siècle, le mathématicien anglais Charles Babbage entreprend la construction d'une machine à calculer appelée *machine aux différences* ou *machine analytique* capable d'effectuer des calculs exacts sans erreurs et de les imprimer. Pour cela, il s'inspire des machines construites par Pascal et Leibniz et a l'idée, en 1834, d'y ajouter un système de cartes perforées comme celui utilisé par Jacquard pour ses métiers à tisser. En 1843, Ada Lovelace publie la traduction d'un article sur la machine de Babbage augmentée de sept notes personnelles. La note G, notamment, contient un algorithme permettant de calculer des nombres particuliers (les nombres de Bernoulli). Cette note contient deux nouveautés majeures : tout d'abord, elle introduit une boucle conditionnelle et, ensuite, son algorithme est écrit, non pas de façon abstraite et générale, mais dans une forme destinée à être appliquée directement à la machine analytique de Babbage. En ce sens, il est considéré comme le premier programme informatique de l'histoire.

À la fin du XIXe siècle, Herman Hollerith, employé au bureau du recensement américain, met au point une machine permettant de lire et de trier des cartes perforées. Ce système sera employé pour le recensement de 1890 qui ne prendra « que » 6 ans, contre 10 ans auparavant. C'est le début de la mécanographie, ancêtre de l'informatique. Par la suite, Hollerith quittera l'administration américaine pour fonder sa propre société qui deviendra, en 1917, IBM.

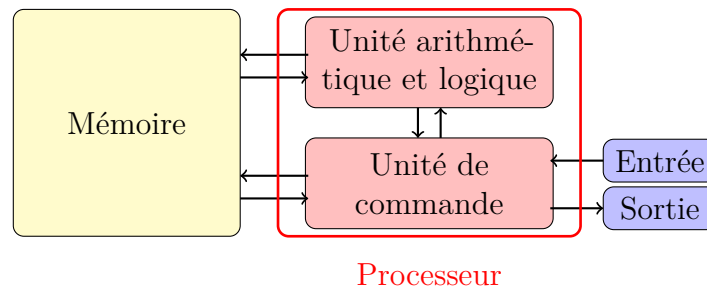
2) Invention de l'ordinateur

Jusqu'au début du XXe siècle, les machines pouvaient exécuter des calculs simples ou des tâches déterminées (comme les métiers à tisser) mais guère plus. En 1936, le mathématicien anglais Alan Turing publie un article dans lequel il définit une machine théorique, appelée depuis *machine de Turing*, qui est plus un concept qu'une machine réelle et qui permet de donner une définition de ce qui est « effectivement calculable » ou plus généralement ce qu'est un algorithme. L'article de Turing comporte deux idées majeures : d'abord, la nécessité de mémoriser des informations et donc de disposer d'une mémoire et, ensuite, le fait d'avoir une machine (théorique) qui peut effectuer n'importe quelle action. En effet, Turing montre l'existence d'*une machine universelle* c'est-à-dire d'une machine de Turing capable de simuler le comportement de n'importe quelle machine de Turing.

Cette vision, bien que purement théorique, est à l'origine du concept d'ordinateur.

En 1945, le mathématicien et physicien américano-Hongrois John von Neumann pose les bases de l'architecture des ordinateurs. Selon ses principes, un ordinateur est « une machine universelle contrôlée par un programme dont les instructions sont codées sous forme numérique (binaire) et enregistrés en mémoire » ⁽¹⁾. L'architecture de von Neumann se caractérise par 4 types d'éléments principaux :

- l'unité arithmétique et logique (UAL) chargée des opérations arithmétiques élémentaires ;
- l'unité de commande chargée du « séquençage » des opérations,
- la mémoire qui contient à la fois les données et le programme chargé d'indiquer à l'unité de commande les opérations à faire sur ces données,
- les unités d'entrée et de sortie qui permettent à la machine de communiquer avec le monde extérieur.



Le modèle de von Neumann

Les premières machines opérationnelles construites selon l'architecture de von Neumann voient le jour au Royaume-Uni avec l'EDSAC (Electronic Delay Storage Automatic Calculator) achevé en mai 1949 à l'université de Cambridge et le Manchester Mark 1 mis au point à la Victoria University de Manchester.

3) Développement des ordinateurs

Les deux premiers ordinateurs à usage commercial furent le Ferranti Mark 1, basé sur le Manchester Mark 1, fabriqué par Ferranti à partir de février 1951 et l'UNIVAC (UNIVersal Automatic Computer) basé sur les travaux d'Eckert et Mauchly à la suite de l'ENIAC et fabriqué par Remington Rand à partir de mars 1951. Ce dernier fut le premier ordinateur à avoir une « large » distribution (1500 exemplaires vendus). Il contenait 5 200 tubes et pesait 13 tonnes.

Par la suite, diverses innovations technologiques vont permettre de rendre les ordinateurs de plus en plus fiables, légers, bon marché et simples d'utilisation pour aboutir à une utilisation généralisée et quotidienne.

- 1949 • Invention du transistor qui remplacera les tubes à vide.
- années 1950 • Ordinateurs de seconde génération.
- 1958 • Invention du circuit intégré.
- années 1960 • Ordinateurs de troisième génération.
- 1969 • Invention du microprocesseur.
- années 1970 • Ordinateurs de quatrième génération.
- 1975 • 2^{de} loi de Moore : la capacité des microprocesseurs double tous les 2 ans.
- années 1980 • Développement des systèmes d'exploitation et des logiciels.
- années 1990 • Développement des réseaux de communication.
- années 2000 • Invention des smartphones.

(1). P. Zanella, Y. Ligier et E. Lazard, *Architecture et technologie des ordinateurs - 5^{ème} édition : Cours et exercices corrigés*, Dunod, 2013, p. 16.

II. — Stockage des données

1) Unité de mémoire

Dans le modèle de von Neumann, qui correspond aujourd'hui encore à la structure globale des ordinateurs, les programmes et les données sont stockés dans un espace mémoire.

Un tel espace est composé de milliards de dispositifs électroniques pouvant être dans 2 états possibles qu'on peut interpréter comme 0 ou 1. Cette unité élémentaire de mémoire à deux états est appelée un *bit* (contraction de *binary digit*). Un agglomérat de 8, 16, 32 ou plus de bits constituent ce qu'on appelle des *cases mémoires*. Le nombre de cases mémoires définit la taille de la mémoire de l'ordinateur.

Une autre unité utilisée est l'octet. Un octet est composé de 8 bits. Comme chaque bit peut prendre deux valeurs (0 ou 1), un octet peut coder $2^8 = 256$ valeurs différentes. Généralement, la taille d'une mémoire est donnée en octet ou dans une puissance de dix d'octets :

unité	kiloctet (ko)	mégaoctet (Mo)	gigaoctet (Go)	teraoctet (To)	pétaoctet (Po)
en octets	10^3	10^6	10^9	10^{12}	10^{15}

2) Support de stockage

L'information a été stockée au cours du temps sur différents types de supports internes ou externes.

La première génération de supports était constituée par les ruban perforés et les cartes perforées. La deuxième génération est constituée par les supports magnétiques. D'abord utilisés sous forme de bandes magnétiques à partir de 1928, ils se sont ensuite miniaturisés et démocratisés sous forme de disques durs à partir des années 1950 puis de cassettes et de disquettes dans les années 1960. La troisième génération est constituée par les supports optiques. Parmi ceux-ci, on trouve le disque compact (CD) développé au début des années 1980, le DVD apparu au milieu des années 1990 et le Blu-Ray commercialisé à partir du milieu des années 2000. La quatrième génération, enfin, est constituée par les mémoires flash. Cette technologie développée à la fin des années 1980 est au début coûteuse et peu efficace pour le stockage. Elle va cependant être améliorée pour devenir le mode standard de nos jours avec les clé USB, les cartes SD et micro SD et les disques SSD.

- 1928 — Bande magnétique : 50 octets par cm.
- 1967 — Disquette 8" : environ 82 ko.
- 1981 — Disquette 5¼" : environ 369 ko.
- 1984 — Disquette 5¼" (2nde génération) : environ 1,2 Mo.
- 1987 — Disquette 3½" : environ 1,47 Mo
- 1990 — CD-ROM : environ 682 Mo
- 1995 — DVD : entre 4,7 et 17 Go.
- 2000 — Mémoire flash : 16 Go.
- 2007 — SSD : 1 To.
- 2013 — SSD : 6 To.
- 2015 — SSD : 10 To.
- 2019 — SSD : 16 To.



bande magnétique



disquettes 8", 5¼" et 3½"



cassette



CDRom



Clé USB



Disque dur SSD

3) Format des données stockées

Un ordinateur peut traiter des données de natures très différentes (texte, image, son, vidéo, logiciel, programme...) une fois que celles-ci ont été numérisées c'est-à-dire transformées en une série de bits. Il existe différentes façons de coder des données qui constituent ce qu'on appelle le format d'un fichier numérique. Le format se reconnaît à l'extension présente dans le nom du fichier c'est-à-dire la chaîne de caractère préfixée par un point qu'on trouve à la fin de ce nom (par exemple, un_texte.txt, un_document_word.docx, un_fichier_son.mp3, un_fichier_video.avi). On peut distinguer deux types de fichiers. D'une part, les fichiers exécutables (comme les logiciels) qui sont écrits dans un langage compréhensible par la machine ou sous une forme interprétable directement par l'ordinateur (c'est-à-dire traduit en langage machine à la volée au fur et à mesure de l'exécution). D'autre part, les fichiers non exécutables qui ne peuvent pas être lus directement par la machine et qui demandent en général de disposer d'un logiciel spécifique pour être lus ou modifiés. Pour une même nature de donnée, il existe souvent plusieurs formats qui nécessitent parfois des logiciels différents pour être ouverts. Le tableau ci-dessous regroupe quelques extensions courantes (mais il y en a beaucoup d'autres).

extension	.txt	.docx	.odt	.jpeg	.png	.mp3	.wav	.avi	.mp4
nature	texte	texte	texte	image	image	son	son	vidéo	vidéo

La taille d'un fichier (c'est-à-dire l'espace mémoire occupé par le fichier) est très variable. Les différents formats utilisent souvent des procédés de numérisation qui permettent d'économiser de la place. Par exemple, le format mp3 n'encodent pas les fréquences que l'oreille humaine ne peut pas entendre ou certains formats n'encodent qu'une seule fois un élément redondant (un même mot utilisé plusieurs fois dans un texte, une même image présente plusieurs fois dans une vidéo...). À titre indicatif, on peut retenir les ordres de grandeurs suivants : un fichier texte est de l'ordre du ko, un fichier son ou un fichier image est de l'ordre du Mo et une vidéo est de l'ordre de quelques centaines de Mo au Go.

4) L'exemple du code ASCII

Le format .txt utilise un code, appelé code ASCII, dans lequel chaque « caractère » (lettres, espaces, ponctuation, retours à la ligne, etc) correspond à un nombre entre 0 et $2^7 - 1 = 127$ dans sa première version et entre 0 et $2^8 - 1 = 255$ dans sa version étendue qui contient les caractères accentués. Par exemple, les lettres majuscules correspondent aux nombres de 65 à 90 et les lettres minuscules aux nombres de 97 à 122. Ces nombres sont ensuite codés en binaire sur un octet (ce qui est possible puisqu'on utilise 2^8 nombres). Par exemple, la lettre **f** correspond au nombre 102 et l'écriture binaire de 102 est 01100110 donc, dans le format .txt, la lettre **f** est codé par l'octet 01100110.

Ce qu'il faut retenir, c'est qu'un caractère est codé par un octet donc pour connaître la taille d'un fichier .txt, il suffit de compter le nombre de caractères. Par exemple, un fichier .txt contenant le texte :

Quelle est la taille de ce texte ?

a une taille de 34 octets car il contient 34 caractères (espaces et ponctuation compris).

III. — L'intelligence artificielle

Si l'idée de construire des machines imitant le comportement humain est ancienne (comme les automates de Vaucanson au milieu du XVIIIe), on peut situer la naissance de l'intelligence artificielle dans les années 1950 avec l'article fondateur d'Alan Turing intitulé *Computing Machinery and Intelligence* et la conférence organisée au Dartmouth College (Hanover, New Hampshire) en 1956 sur le thème des machines pensantes durant laquelle le chercheur américain John McCarthy introduit pour la première fois le terme *intelligence artificielle*.

1) Définition

Il n'est pas simple de clairement définir ce qu'on appelle une machine « intelligente ». Est-ce une machine capable d'apprendre ? d'utiliser des connaissances ? d'interagir avec son environnement ? de planifier une succession de tâches ?

Le notion d'**intelligence artificielle** (abrégé IA) possède des contours encore assez flous. On peut cependant adopter la définition suivante, due au chercheur français Yann Le Cun :

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

Ainsi, l'IA a pour but de rendre des machines capables de reproduire des activités humaines, que ces activités soient de l'ordre de la compréhension, de la perception ou de la décision. On peut citer par exemple :

- la reconnaissance d'objets, d'animaux ou de personnes sur une photo ;
- l'apprentissage de jeux (comme les échecs) ;
- le pilotage de voiture ;
- la traduction de texte ;
- l'établissement de diagnostic médical.

Pour atteindre ce but, une des techniques utilisées (mais à laquelle ne se réduit pas pour autant l'IA) est l'apprentissage machine.

2) Apprentissage machine

L'**apprentissage machine** (appelé aussi apprentissage artificiel ou apprentissage automatique) est un domaine de recherche commun à l'IA et aux statistiques. Il consiste à élaborer des programmes dont le comportement peut évoluer en fonction de données dites *données d'entraînement*.

Le principe général est le suivant : on fournit à la machine un très grand nombre de données à partir desquelles la machine s'entraîne (phase d'apprentissage ou d'entraînement) afin de déterminer le comportement qu'elle adoptera ultérieurement sur de nouvelles données (phase d'inférence).

Imaginons, par exemple, qu'on souhaite apprendre à une machine à reconnaître des images de chats. Plutôt que d'essayer d'écrire des algorithmes complexes permettant d'identifier différentes caractéristiques de l'animal, on fournit à la machine un grand nombre d'images dont certaines sont des images de chat et d'autres non. La machine va alors déterminer des paramètres qui vont permettre de distinguer les photos contenant des chats des autres photos : c'est la phase d'apprentissage. À la fin de cette phase, la machine a développé ses propres critères de choix et elle peut les appliquer à des nouvelles photos : c'est la phase d'inférence. Notons ici une grande différence dans la performance de l'apprentissage machine par rapport à l'apprentissage humain : là où moins de 10 photos suffisent à un enfant pour savoir reconnaître un chat, il en faut plusieurs milliers pour une machine !

En utilisant les données d'entraînement, la machine va construire une fonction d'apprentissage destinée à répondre à une situation donnée. Cette fonction peut être une fonction de décision (comme dans le cas des photos de chats) mais elle peut aussi consister à déterminer un (ou plusieurs) nombres.

On distingue trois grands type d'apprentissages machine en fonction des modalités de la phase d'entraînement.

- **l'apprentissage supervisé** : dans ce type d'apprentissage, la machine s'entraîne sur des données qui ont été au préalable étiquetées par l'homme. Si on reprend l'exemple précédent, on fournira à la machine des images en indiquant celles qui représentent des chats et celles qui n'en représentent pas.

Parmi les techniques d'apprentissage supervisé, on peut signaler l'apprentissage profond qui, schématiquement, consiste en une succession de modules tels que les résultats produits par les uns sont utilisés comme données d'entrée par les autres. Ce principe s'inspire directement de l'architecture du cerveau humain sous la forme de réseaux de neurones artificiels.

- **l'apprentissage non supervisé** : dans ce type d'apprentissage, on fournit à la machine des données non étiquetées et on la laisse repérer des régularités, des proximités, des corrélations pour construire elle-même la fonction d'apprentissage.
- **l'apprentissage par renforcement** : dans ce type d'apprentissage, la machine va faire des essais et bénéficier d'un système de récompenses/punitions en fonction du succès de ses actions. C'est, par exemple, le cas des programmes qui apprennent à jouer à des jeux du type jeu d'échecs ou jeu de go. Dans ce cas, la machine va jouer (contre-elle même par exemple) et, grâce au système de récompenses, créer une fonction d'apprentissage qui lui permettra de trouver l'action optimale en fonction de la situation de jeu. Ainsi, ici, la machine crée ses propres données d'entrée en jouant des parties d'entraînement.

Les différentes méthodes s'appuient en général sur des outils mathématiques pour construire la fonction d'apprentissage.

On peut en donner quelques exemples dans des cas particulièrement simples.

- **Utilisation de courbes d'ajustement** : lorsqu'on dispose de données reliant deux paramètres (population en fonction du temps, prix de l'immobilier en fonction du temps, risque de développer une maladie en fonction de la présence d'un certain marqueur, etc.), on peut essayer de trouver une fonction dont la courbe s'ajuste le mieux possible aux données. Cette courbe peut être une droite (comme on l'a vu dans le thème 3.4 pour le modèle linéaire) mais peut prendre d'autres formes. Les données d'entraînement permettent dans ce cas de déterminer un modèle qui sera d'autant plus pertinent que le nombre de données sera grand.

- **Méthode du (ou des) plus proche(s) voisin(s)** : en reprenant la situation précédente, dans le cas où on souhaite classer les données en deux groupes A et B, on peut placer les données d'entraînement sous la forme de points dans un repère et, lorsqu'on doit traiter une nouvelle donnée, on place le point correspondant dans un repère et on cherche son plus proche voisin, c'est-à-dire le point issu des données d'entraînement le plus proche du nouveau point. Si ce plus proche voisin fait partie du groupe A, on classe la nouvelle donnée dans le groupe A et, sinon, on la classe dans le groupe B.

Une autre possibilité est de considérer non pas seulement le plus proche voisin mais les k plus proches voisins (où k est un entier impair). Si, parmi ces k plus proches voisins, une majorité appartient au groupe A, on classe la nouvelle donnée dans le groupe A et, sinon, on la classe dans le groupe B.

- **Inférence bayésienne** : il s'agit d'une méthode de détermination des causes à partir des conséquences basées sur des calculs de probabilités. On détaillera cette technique dans la paragraphe suivant.

IV. — L'inférence bayésienne

1) Principe

L'inférence bayésienne est une méthode de calcul permettant de déterminer les probabilités des causes à partir de probabilités de leurs effets. L'adjectif « bayésienne » provient du nom d'un révérend et mathématicien anglais du XVIII^e siècle, Thomas Bayes, à qui l'on doit notamment un théorème très utilisé en probabilité.

L'inférence bayésienne est utilisée en apprentissage automatique, notamment dans le cas de prise de décision : une personne ayant un test positif pour une maladie est-elle effectivement malade ? un courrier électronique ayant un contenu suspect doit-il être considéré comme un spam ?

L'inférence bayésienne va permettre, à partir d'un grand nombre de données, de fournir une réponse probabiliste à ces questions.

2) Exemple d'un test de dépistage

Pour un test de dépistage, on définit les deux caractéristiques suivantes :

- sa **sensibilité** qui représente la probabilité qu'une personne malade effectuant ce test ait un résultat positif ;
- sa **spécificité** qui représente la probabilité qu'une personne non malade effectuant ce test ait un résultat négatif.

Cependant, lorsqu'une personne a un test positif, ce qui importe est de savoir si elle est vraiment malade.

On va voir que cela dépend essentiellement de la **prévalence** de la maladie c'est-à-dire la proportion de personnes atteintes de la maladie dans la population.

Considérons, par exemple, une population de 100 000 habitants et une maladie dont la prévalence est 0,1%. Supposons qu'on dispose d'un test dont la sensibilité est 95% et la spécificité est 98%. Ainsi, dans cette population, il y a :

-
-
-
-

On peut rassembler ces valeurs dans le tableau de contingence suivant :

	Test positif	Test négatif	Total
Personnes malades			
Personnes non malades			
Total			100 000

Ainsi,

- la **valeur prédictive positive**, c'est-à-dire la probabilité qu'une personne ayant un test positif soit effectivement malade est égale à
- la **valeur prédictive négative**, c'est-à-dire la probabilité qu'une personne ayant un test négatif soit effectivement non malade est égale à

Ainsi, on constate que, même si la sensibilité et la spécificité sont très bonnes, la valeur prédictive positive est faible c'est-à-dire une personne testée positive a peu de chance d'être effectivement malade.

Ce paradoxe apparent s'explique par le fait que les valeurs prédictives dépendent en fait grandement de la prévalence de la maladie.

La formule de Bayes permet de montrer que, si un test a une sensibilité Se et une spécificité Sp alors la valeur prédictive positive de ce test est

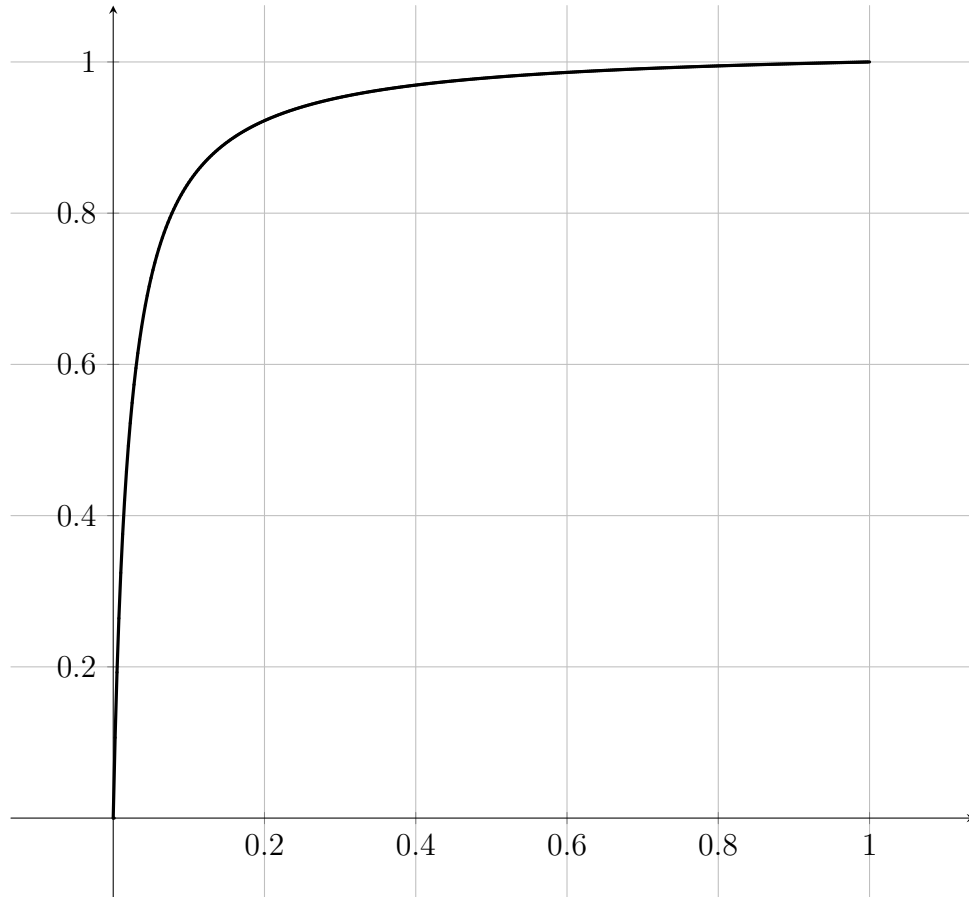
$$V = \frac{p \times Se}{p \times Se + (1 - p)(1 - Sp)}$$

où p désigne la prévalence de la maladie dans la population.

Ainsi, si on considère le test précédent pour lequel $Se = 0,95$ et $Sp = 0,98$, on obtient une valeur prédictive positive en fonction de p égale à

$$V(p) = \frac{p \times 0,95}{p \times 0,95 + (1 - p) \times 0,02} = \frac{0,95p}{0,93p + 0,02}$$

On a tracé ci-dessous la courbe représentative de la fonction V .



On constate que pour des prévalences faibles, la valeur prédictive positive est également faible, et ceci pour une sensibilité et une spécificité constantes.

3) Détection de spams

Un des premiers programmes de filtrage bayésien du courrier électronique a été le programme iFile conçu par Jason Rennie et publié en 1996. Le principe, analogue à celui du diagnostic médical, repose sur le fait que les mots du dictionnaire ont des probabilités différentes d'apparaître dans les spams et dans les courriers légitimes. Le filtre de détection des spams ne connaît pas à l'avance les probabilités d'apparition de ces mots, c'est pourquoi il lui faut une phase d'apprentissage pour les évaluer.

L'apprentissage se fait à partir de l'observation du comportement des utilisateurs, qui doivent indiquer manuellement si un message est un spam ou non. Le filtre ajuste les probabilités de rencontrer un mot M donné dans un spam ou dans un courrier légitime à l'aide des messages d'entraînement. Ensuite, en utilisant la formule de Bayes, la machine calcule la probabilité que le message soit un spam sachant qu'il contient le mot M . Cette probabilité est enfin comparée à un seuil : si elle est supérieure au seuil, le filtre classera ce message dans les spams.

Stockage des données

Exercice 1. Quel type de croissance est décrit par la seconde loi de Moore ?

Exercice 2. Avant d’être mesurées par des puissances de 10 d’octets (ko, Mo, Go, etc), les données informatiques étaient mesurées en puissances de 2 d’octets. Pendant longtemps, un kilooctet a désigné 2^{10} octets, un mégaoctet a désigné 2^{20} octets, un gigaoctet a désigné 2^{30} octets et ainsi de suite. Cependant, cette tradition propre au domaine de l’informatique entraine en conflit avec les normes internationales (selon lesquelles 1 kilo correspond à 10^3 , un méga à 10^6 , etc). Ainsi, en 1988, la Commission électrotechnique internationale a normalisé les unités en introduisant des préfixes spécifiques pour les puissances de 2. Ainsi, sont apparus les termes *kilo binaire*, *méga binaire*, *giga binaire*, *téra binaire* et *péta binaire* abrégé en kibi, mébi, gibi, tébi et pébi. On a alors le tableau suivant :

unité	kibioctet (Kio)	mébioctet (Mio)	gibioctet (Gio)	tébioctet (Tio)	pébioctet (Pio)
en octets	2^{10}	2^{20}	2^{30}	2^{40}	2^{50}

1. Les premières disquettes 3 1/2 construites par Sony avaient une capacité de 400 Kio. Déterminer la capacité en ko d’une telle disquette.
2. Un disque dur S-ATA Hitachi de fin 2005 avait une capacité de stockage de 76,688 Gio. Convertir cette capacité en Go.

Exercice 3.

1. Quelle la taille en octet d’un fichier texte codé en ASCII et contenant le texte suivant ?

Je dois déterminer la taille de ce texte.
 Pour cela, je ne dois pas oublier les espaces et la ponctuation,
 ni les retours à la ligne.

2. Un fichier texte codé en ASCII compte 12 lignes. Chaque ligne compte 30 caractères (espaces et ponctuation compris). Quelle est la taille en octet de ce fichier ?
3. Un fichier texte codé en ASCII a un taille de 236 ko. Déterminer le nombre maximum de caractères qu’il peut contenir.

Exercice 4. Dans un dossier, on trouve les fichiers suivants :

fichier1.jpg fichier2.txt fichier3.mov fichier4.exe fichier5.png
 fichier6.mp4 fichier7.doc fichier8.avi fichier9.wav fichier10.mp3

Regrouper ces fichiers en 5 catégories : fichiers texte, fichiers image, fichier son, fichiers vidéo et fichiers exécutables.

Exercice 5. On dispose 3 fichiers `fichier1`, `fichier2` et `fichier3`. On sait que la taille de `fichier1` est 840 Mo, la taille de `fichier2` est 53 Mo et la taille de `fichier3` est 14 ko. On sait également que les extensions de ces fichiers sont `.txt`, `.wav` et `.avi`.

En se référant aux ordres de grandeurs standards, déterminer l’extension de chaque fichier.

Exercice 6. (Source : <http://images.math.cnrs.fr/Le-traitement-numerique-des-images.html>)

Lorsqu'on numérise une image en niveau de gris, on stocke cette image sous la forme d'un tableau constitué de petits carrés appelés pixels. À chaque pixel, on va associer un nombre entre 0 et 255 correspondant à un certain niveau de gris comme sur l'image suivante :



Image A

Chaque nombre est ensuite codé sur un octet comme dans le cas du code ASCII.

1. L'image A a une résolution de 240 par 240, c'est-à-dire qu'elle correspond à un tableau ayant 240 lignes et 240 colonnes. Combien de pixels constituent cette image ? Déterminer sa taille en ko.
2. Pour gagner de la place, on peut soit diminuer le nombre de pixels soit diminuer le nombre de niveaux de gris.



Image B



Image C

- a. Déterminer la taille de l'image B sachant qu'elle a été obtenue à partir l'image A en ne conservant qu'une ligne sur deux et qu'une colonne sur deux.
- b. Déterminer la taille de l'image C sachant qu'elle a été obtenue à partir l'image A en ne conservant que 16 niveaux de gris au lieu de 256.

Exercice 7.

1. Écrire une fonction Python `en_binaire(N)` qui prend en argument un entier `N` compris en 0 et 255 et qui renvoie l'octet (sous forme d'une chaîne de caractères) qui code cet entier en écriture binaire.
2. Écrire une fonction Python `en_decimal(O)` qui prend en argument un octet `O` (sous forme d'une chaîne de caractères) et qui renvoie l'entier codé par `O` en binaire.
3. Déterminer le code ASCII binaire correspondant au mot `Binaire`.
4. À quel mot correspond le code ASCII binaire suivant :

01000100 01100101 01100011 01101001 01101101 01100001 01101100

Stockage des données – Corrigés

Exercice 1. La seconde loi de Moore postule que la capacité des microprocesseurs (en équivalents de nombres de transistors) double tous les 2 ans. Il s'agit donc d'une croissance exponentielle.

Exercice 2.

1. Comme 400 Kio correspond à $400 \times 2^{10} = 409\,600$ octets, la capacité en ko des premières disquettes de Sony était d'environ 410 ko.
2. Comme 76,688 Gio correspond à $76,688 \times 2^{30} \approx 82,3 \times 10^9$ octets, un disque dur S-ATA Hitachi de fin 2005 avait une capacité de stockage d'environ 82,3 Go.

Exercice 3.

1. Le texte compte 104 lettres, 4 caractères de ponctuation, 23 espaces et 2 retours à la ligne donc la taille du fichier texte codé en ASCII est $104 + 4 + 23 + 2 = 133$ octets.
(Dans la pratique, on obtient un fichier de 135 octets car il semble que les caractères accentués soient codés sur 2 octets et non pas 1 octet comme les autres caractères.)
2. Le fichier compte $12 \times 30 = 360$ caractères ainsi que 11 retours à la ligne donc sa taille est $360 + 12 = 372$ octets.
3. Chaque caractère occupant 1 octet, un fichier de 236 ko contient au maximum 236 caractères.

Exercice 4. On peut classer les fichiers de la façon suivante :

- fichiers texte : `fichier2.txt` et `fichier7.doc`
- fichiers image : `fichier1.jpg` et `fichier5.png`
- fichiers son : `fichier9.wav` et `fichier10.mp3`
- fichiers vidéo : `fichier3.mov`, `fichier6.mp4` et `fichier8.avi`
- fichier exécutable : `fichier4.exe`

Exercice 5. En se référant aux ordres de grandeurs standards, `fichier1` est une vidéo donc son extension est `.avi`, `fichier2` est un fichier son donc son extension est `.wav` et `fichier3` est un fichier texte donc son extension est `.txt`.

Exercice 6.

1. L'image contient $240 \times 240 = 57\,600$ pixels donc, comme chaque pixel est codé sur un octet, la taille de l'image est 57,6 ko.
2. a. En ne conservant qu'une ligne sur deux et qu'une colonne sur deux, on obtient une image comptant $120 \times 120 = 14\,400$ pixels et sa taille est donc 14,4 ko.
b. En ne conservant que 16 niveaux de gris, on peut coder chaque niveau de gris sur 3 bits (car $2^3 = 8$) et donc chaque pixel nécessite 3 bits de mémoire. Ainsi, la taille de l'image est $57\,600 \times 3 = 172\,000$ bits c'est-à-dire $\frac{172\,000}{8} = 21\,600$ octets soit finalement 21,6 ko.

Exercice 7.

1.

```
def en_binaire(N):
    O=' '
    for i in range(8):
        b=str(N//2**(7-i))
        O=O+b
        N=N%2**(7-i)
    return(O)
```

2.

```
def en_decimal(O):
    N=0
    for i in range(8):
        N=N+int(O[7-i])*2**i
    return(N)
```

3. On utilise la table ASCII suivante :

Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	@	96	60	`
1	1	[START OF HEADING]	33	21	!	65	41	A	97	61	a
2	2	[START OF TEXT]	34	22	"	66	42	B	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	C	99	63	c
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	e
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	'	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	H	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49	I	105	69	i
10	A	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	B	[VERTICAL TAB]	43	2B	+	75	4B	K	107	6B	k
12	C	[FORM FEED]	44	2C	,	76	4C	L	108	6C	l
13	D	[CARRIAGE RETURN]	45	2D	-	77	4D	M	109	6D	m
14	E	[SHIFT OUT]	46	2E	.	78	4E	N	110	6E	n
15	F	[SHIFT IN]	47	2F	/	79	4F	O	111	6F	o
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	P	112	70	p
17	11	[DEVICE CONTROL 1]	49	31	1	81	51	Q	113	71	q
18	12	[DEVICE CONTROL 2]	50	32	2	82	52	R	114	72	r
19	13	[DEVICE CONTROL 3]	51	33	3	83	53	S	115	73	s
20	14	[DEVICE CONTROL 4]	52	34	4	84	54	T	116	74	t
21	15	[NEGATIVE ACKNOWLEDGE]	53	35	5	85	55	U	117	75	u
22	16	[SYNCHRONOUS IDLE]	54	36	6	86	56	V	118	76	v
23	17	[ENG OF TRANS. BLOCK]	55	37	7	87	57	W	119	77	w
24	18	[CANCEL]	56	38	8	88	58	X	120	78	x
25	19	[END OF MEDIUM]	57	39	9	89	59	Y	121	79	y
26	1A	[SUBSTITUTE]	58	3A	:	90	5A	Z	122	7A	z
27	1B	[ESCAPE]	59	3B	;	91	5B	[123	7B	{
28	1C	[FILE SEPARATOR]	60	3C	<	92	5C	\	124	7C	
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D]	125	7D	}

Ainsi, Binaire se traduit en décimal 66 105 110 97 105 114 101 et, grâce à la fonction en_binaire, on en déduit que le code ASCII du mot Binaire est

01000010 01101001 01101110 01100001 01101001 01110010 01100101

4. Grâce à la fonction en_decimal, on trouve que la suite d'octets

01000100 01100101 01100011 01101001 01101101 01100001 01101100

correspond à la suite d'entiers 68 101 99 105 109 97 108 et donc, en utilisant la table précédente, on conclut que le mot correspondant est Decimal.

Intelligence artificielle

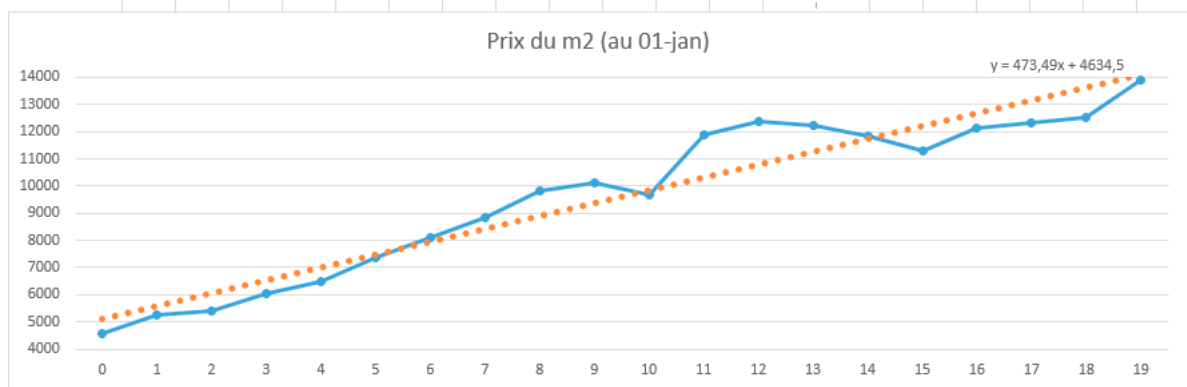
Exercice 1. Dans le cours, on a mentionné que les logiciels d'apprentissage de jeux de type jeu d'échecs ou jeu de go utilisent un apprentissage par renforcement en jouant des parties d'entraînement.

Comment devrait-on s'y prendre pour entraîner un logiciel d'apprentissage de jeux en utilisant un apprentissage supervisé ?

Exercice 2. Le tableau suivant donne le prix moyen du m² dans le 6^e arrondissement de Paris au 1^{er} janvier entre 2000 et 2019. Ces données sont représentées sur le graphique en dessous et on a également tracé en pointillés la droite d'ajustement de ces données.

Evolution du prix du mètre carré dans le 6^e arrondissement de Paris (année 0 en 2000)

Année (x _i)	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Prix du m ² y _i (au 01-jan)	4562	5270	5400	6020	6460	7360	8090	8840	9830	10100	9690	11870	12400	12250	11820	11280	12150	12320	12530	13880



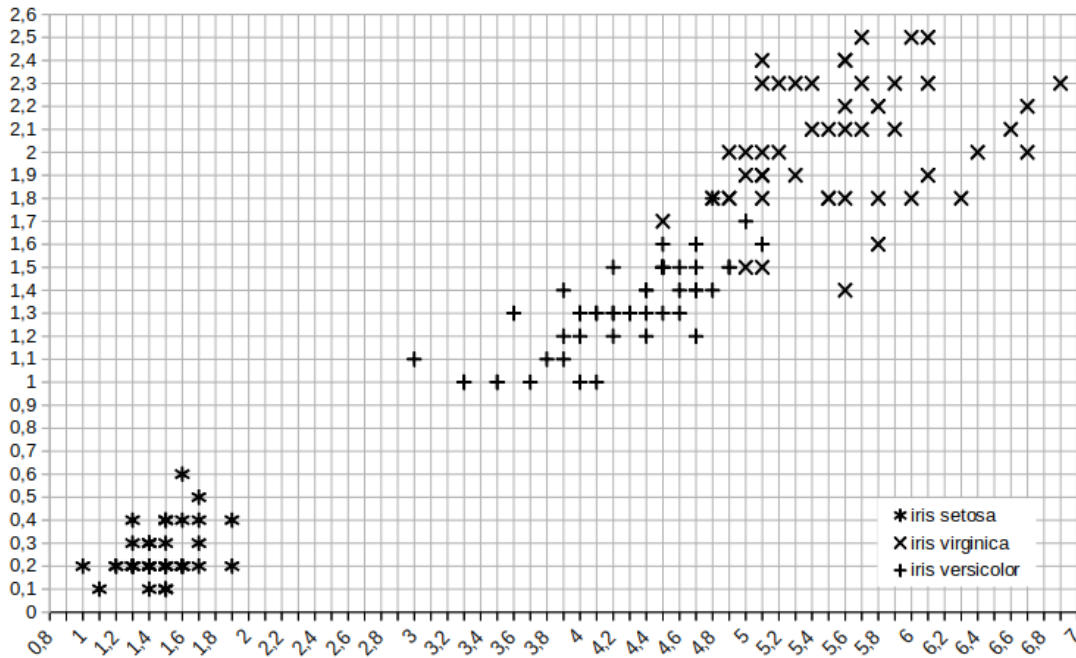
Donner une estimation du prix du mètre carré dans le 6^e arrondissement au 1^{er} janvier 2025.

Exercice 3. On se propose de répartir en deux catégories (Maligne, Bénigne) des tumeurs dont on connaît un certain nombre de caractéristiques. Dans la réalité, un tel diagnostic automatique est effectué à partir de données d'apprentissage portant sur une cinquantaine de caractéristiques mesurées sur un échantillon d'un millier de tumeurs déjà étiquetées « Maligne » ou « Bénigne ». Dans un souci de simplification, on se restreint ici à 2 caractéristiques (diamètre et concavité) mesurées sur un panel de 10 patientes.

Diamètre moyen (en mm)	Concavité moyenne	Catégorie
13,2	8,3	Bénigne
18,7	19,7	Bénigne
8,2	15,9	Maligne
13,2	9	Bénigne
13,5	4,8	Maligne
11,8	1,7	Maligne
13,6	1,9	Maligne
12	2	Maligne
18,2	17,7	Bénigne
12	6,6	Maligne

1. À quel type d'apprentissage a-t-on affaire ici ?
2. Placer les points correspondant aux données du tableau dans un repère (on mettra le diamètre en abscisse et la concavité en ordonnées).
3. En utilisant la méthode du plus proche voisin, une tumeur de 10 mm et ayant une concavité de 12 doit-elle être considérée comme bénigne ou maligne ?

Exercice 4. En 1936, Edgar Anderson a collecté des données sur 3 espèces d'iris : iris setosa, iris virginica et iris versicolor. Pour chaque espèce, Anderson a mesuré (en cm) différents paramètres. Sur le graphique ci-dessous, on a représenté la longueur (en abscisse) et la largeur (en ordonnée) mesurées lors de ces relevés.



Déterminer à quelle espèce appartient un iris dont les pétales mesure 2,4 cm de long et 0,8 cm de large :

1. en utilisant le méthode du plus proche voisin ;
2. en utilisant la méthode des 5 plus proches voisin.

Exercice 5. Chercher des arguments et des situations montrant l'importance du choix des données d'entraînement pour une IA.

Exercice 6. Chercher des arguments et des situations montrant les problèmes juridiques et éthiques posés par l'IA.

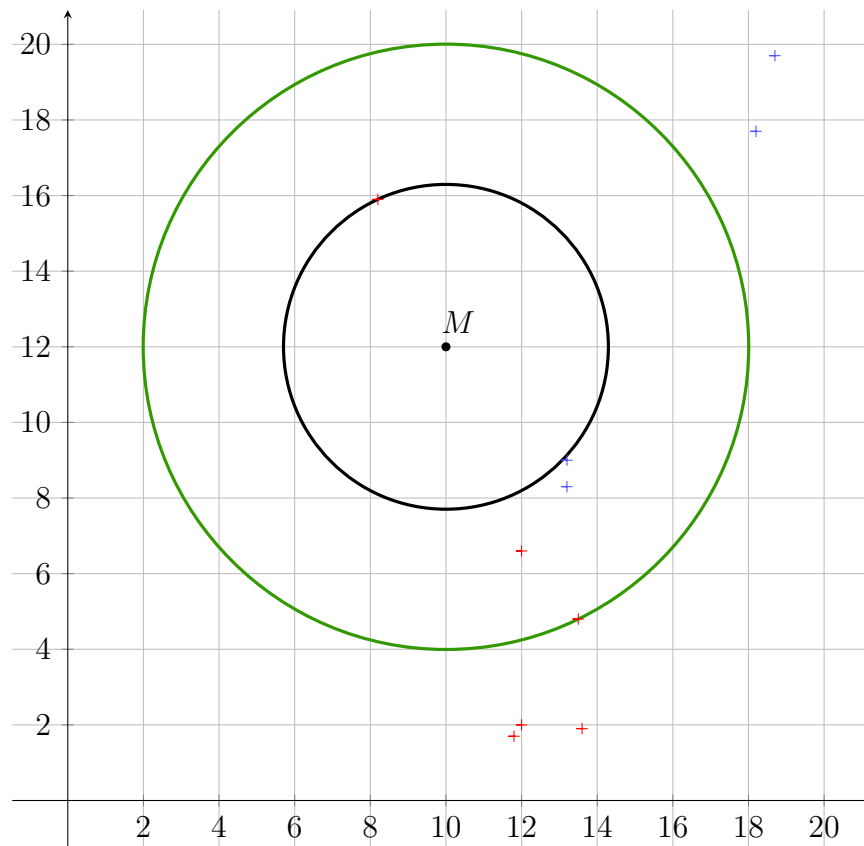
Intelligence artificielle – Corrigés

Exercice 1. Pour entraîner un logiciel d'apprentissage de jeux en utilisant un apprentissage supervisé, il faudrait lui fournir un grand nombre d'exemples de parties jouées.

Exercice 2. La droite d'ajustement affine a pour équation $y = 473,49x + 4634,5$ où x représente le nombre d'années écoulées depuis 2000. Ainsi, on peut estimer le prix du mètre carré dans le 6^e arrondissement au 1^{er} janvier 2025 à $473,49 \times 25 + 4634,5 \approx 16\,470$ euros.

Exercice 3.

1. Étant donné qu'on précise la catégorie des données (bénigne ou maligne), il s'agit d'un apprentissage supervisé.
2. On a représenté les tumeurs bénignes en bleu et tumeurs malignes en rouge.

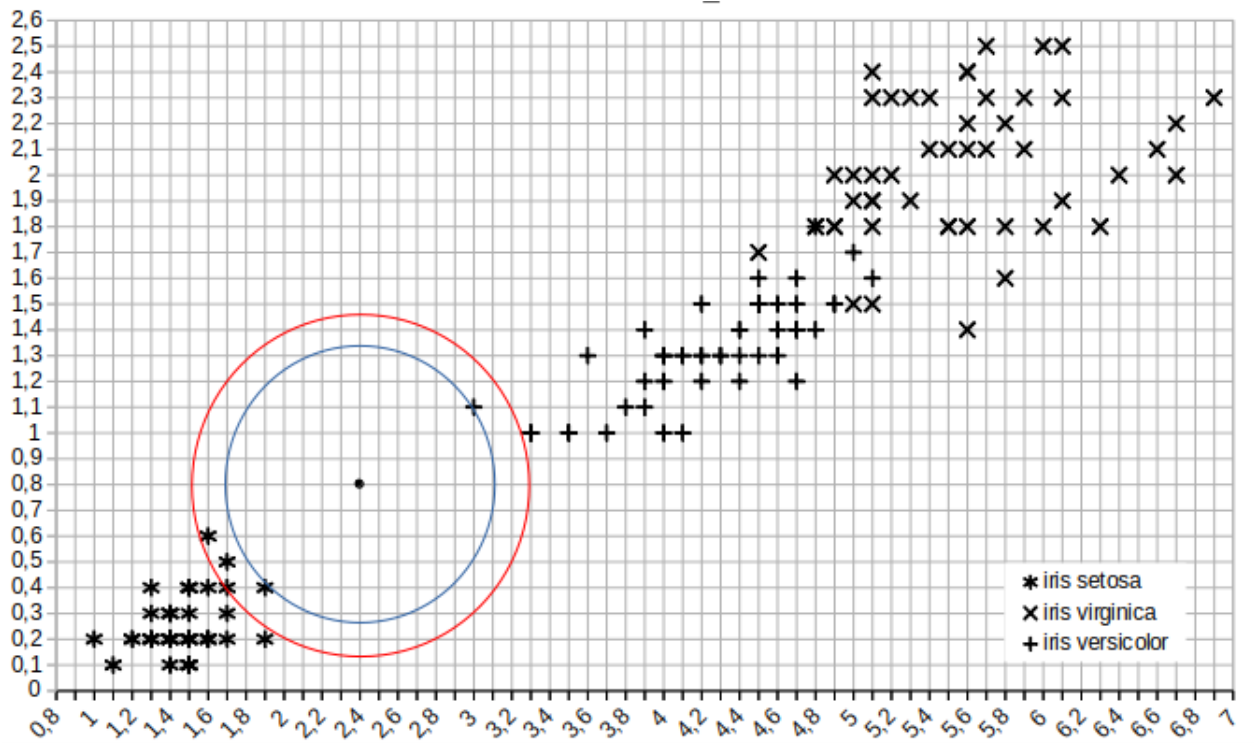


3. La méthode du plus proche voisin indique que la tumeur doit être considérée comme maligne. Cependant, on voit que le point de coordonnées (13,2; 9) est très proche du cercle noir, et quasiment à la même distance de M que le point de coordonnées (8,2; 15,9). Cependant, en utilisant la méthode des 5 plus proches voisins (cercle vert), on retrouve le même résultat.

Remarque. Dans cet exercice, le nombre de données est en fait insuffisant pour obtenir une réponse vraiment pertinente.

Exercice 4. Avec la méthode du plus proche voisin (cercle bleu), on obtient, comme dans l'exercice précédent, un résultat peu concluant car deux points sont quasiment équidistants du point de coordonnées (2,4; 0,8).

Avec la méthode des 5 plus proches voisins (cercle rouge), on obtient un résultat plus net puisque 4 points correspondent à un iris setosa contre 1 seul pour un iris versicolor. Avec cette méthode, on peut donc penser que l'iris considéré est un iris setosa.



Exercice 5. Pour que les données d'entraînement soient efficaces, il est nécessaire qu'elles soient sûres, en quantité suffisante et représentatives.

Par exemple, si les données sont mal mesurées ou périmées, on risque d'obtenir un biais lors de l'apprentissage. L'utilisation de données issues de choix humains peut également reproduire des biais existants. Par exemple, si on entraîne une machine pour évaluer des CV à partir des sélections précédentes faites par des humains, on risque de reproduire des discriminations sexuelles ou ethniques (consientes ou inconsientes).

Utiliser trop peu de données ne permettra pas un calibrage efficient des paramètres de l'algorithme et risque de donner par la suite des mauvais résultats.

Le représentativité des données est un point essentiel. Si on considère l'exemple des photos de chats vu en cours, si toutes les photos de chat représentent des chats noirs, la machine risque d'intégrer la couleur noire parmi les paramètres déterminant et ainsi ne pas reconnaître des chats roux ou blancs.

Exercice 6. L'utilisation de l'IA pose de nombreux problèmes éthiques. Comme on l'a vu, la phase d'entraînement nécessite souvent de disposer de nombreuses données et la récolte et le stockage de ces données peuvent poser des problèmes. De plus, comme toute technologie, l'IA peut être utilisée à des fins critiquables : surveillance des individus, recherche de personnes vulnérables ou influençables, diffusion de fausses nouvelles sur les réseaux sociaux ou uniformisation des contenus pour « satisfaire » l'utilisateur. De plus, les IA peuvent être présentées comme infaillibles et il y a un risque d'instrumentalisation de celles-ci à des fins commerciales, idéologiques ou politiques. Cependant, il ne faut pas oublier que les erreurs et les biais (volontaires ou involontaires) dans la conception, le calibrage et l'utilisation des IA existent et montrent qu'on ne peut pas les considérer comme étant sans défaut.

Sur le plan juridique, il y a également de réels problèmes de responsabilité. Par exemple, si une IA se trompe lors d'un diagnostic médical, qui est responsable de l'erreur : les personnes ayant récoltées les données d'entraînement ? le concepteur du programme ? le médecin qui l'utilise ?

Le même type de questions se pose pour les voitures autonomes. En cas d'accident grave, qui doit être tenu pour responsable ?

Inférence bayésienne

Exercice 1. Le virus de l’immunodéficience humaine (VIH) est responsable du SIDA. Il existe aujourd’hui des tests rapides appelés « Tests rapides d’orientation diagnostique » (TROD) qui ont l’avantage de pouvoir être réalisés à partir d’un échantillon de salive ou d’une goutte de sang prélevée sur le bout du doigt.

1. On a testé avec les deux types de TROD (salivaire et sanguin) deux populations : une première composée de 10 000 personnes infectées par le VIH et une seconde composée de 100 000 personnes non infectées. On a obtenu les résultats suivants.

	Personnes infectées	Personnes non infectées
Tests salivaires positifs	9 803	260
Tests sanguins positifs	9 968	90

Calculer la sensibilité et la spécificité de chacun des deux tests.

2. a. En 2017, la population mondiale exposée au VIH était estimée à 6 milliards et, dans cette population, le nombre de personnes infectées par le VIH était estimé à 37 millions.

Calculer la valeur prédictive positive de chacun des deux tests pour la population mondiale exposée.

- b. En 2017, la population française exposée était estimée à 50 millions et, dans cette population, le nombre de personnes infectées par le VIH était estimé à 150 000.

Calculer la valeur prédictive positive de chacun des deux tests pour la population française exposée.

- c. En 2017, la population sud-africaine exposée était estimée à 35 millions et, dans cette population, le nombre de personnes infectées par le VIH était estimé à 7 millions.

Calculer la valeur prédictive positive de chacun des deux tests pour la population sud-africaine exposée.

- d. Que mettent en évidence les trois exemples précédents ?

Exercice 2. Parmi les femmes de 40 ans ayant effectué une mammographie, 1% a un cancer du sein. À la suite de mammographies sur un échantillon, on a établi que :

- pour 82% des femmes ayant un cancer du sein, la mammographie détecte une anomalie ;
- pour 9% des femmes n’ayant pas de cancer du sein, la mammographie détecte une anomalie.

On suppose que 10 000 femmes de 40 ans ont effectué une mammographie.

1. Déterminer la sensibilité et la spécificité d’une mammographie.
2. Compléter le tableau suivant.

	Anomalie détectée	Pas d’anomalie détectée	Total
Personnes malades			
Personnes non malades			
Total			10 000

3. Une femme de 40 ans a effectué une mammographie qui a permis de détecter une anomalie. Quelle est la probabilité qu'elle soit atteinte d'un cancer du sein ?
4. Calculer les valeurs prédictives positive et négative d'une mammographie chez les femmes de 40 ans.

Exercice 3. On suppose que la somme de la sensibilité et de la spécificité d'un test est égale à 1 (c'est-à-dire, avec les notations du cours, $Se + Sp = 1$).

Expliquer pourquoi le test est inutile dans ce cas.

Les deux exercices suivants sont destinés aux élèves suivant la spécialité mathématiques ou l'option mathématiques complémentaires.

Exercice 4. Trois maladies virales peuvent être transmises par les moustiques : dengue, chikungunya et zika. Elles provoquent des symptômes qui peuvent être assez proches et il est donc parfois difficile de les différencier directement. On s'intéresse à la mise en place d'une aide statistique au diagnostic. Pour cela, on va s'appuyer sur des données obtenues chez des personnes dont le diagnostic a pu être certifié par des examens biologiques. Pour simplifier, on supposera que ces caractères apparaissent indépendamment chez les personnes infectées.

Symptômes	Dengue	Chikungunya	Zika
Fièvre	95%	75%	75%
Courbatures	75%	95%	50%
Douleur oculaire	50%	25%	50%
Déficit en globules blancs	50%	50%	25%
Hémorragies	25%	5%	5%

À partir de ces données, on veut déterminer les probabilités de chaque maladie selon les symptômes présentés et dans des conditions différentes.

1. On suppose qu'une personne malade revient d'un pays dans lequel aucune de ces maladies n'est épidémique. On considère donc a priori que les trois maladies sont équiprobables.
 - a. Quelles sont les probabilités de chaque maladie si cette personne présente à la fois de la fièvre, pas de courbatures et des douleurs oculaires ?
 - b. Quel est le diagnostic le plus probable dans ce cas ?
2. On suppose qu'une personne malade revient d'un pays dans lequel sévit une épidémie de Zika. On suppose qu'il y a 80% de chances qu'elle ait été infectée par Zika et 10% par chacune des deux autres maladies.
 - a. Quelles sont les probabilités de chaque maladie si cette personne présente à la fois de la fièvre, pas de courbatures et des douleurs oculaires ?
 - b. Quel est le diagnostic le plus probable dans ce cas ?

Exercice 5. À l'aide de la formule des probabilités totales, démontrer la formule de Bayes.

Inférence bayésienne – Corrigés

Exercice 1.

1. Pour le TROD salivaire, la sensibilité est $\frac{9\,803}{10\,000} = 0,9803$ et la spécificité $\frac{100\,000 - 260}{100\,000} = 0,9974$.

Pour le TROD sanguin, la sensibilité est $\frac{9\,968}{10\,000} \approx 0,9968$ et la spécificité $\frac{100\,000 - 90}{100\,000} = 0,9991$.

2. a. Sur l'ensemble de la population mondiale en 2017, la proportion de personnes infectées par le VIH est $p = \frac{37\,000\,000}{6\,000\,000\,000} = \frac{37}{6000}$ donc la valeur prédictive du TROD salivaire est

$$\frac{p \times 0,9803}{p \times 0,9803 + (1 - p) \times (1 - 0,9974)} \approx 0,7$$

et la valeur prédictive du TROD sanguin est

$$\frac{0,006 \times 0,997}{0,006 \times 0,997 + (1 - 0,006) \times (1 - 0,999)} \approx 0,873.$$

b. Sur l'ensemble de la population française en 2017, la proportion de personnes infectées par le VIH est $p = \frac{150\,000}{50\,000\,000} = 0,003$ donc la valeur prédictive du TROD salivaire est

$$\frac{0,003 \times 0,9803}{0,003 \times 0,9803 + (1 - 0,003) \times (1 - 0,9968)} \approx 0,532$$

et la valeur prédictive du TROD sanguin est

$$\frac{0,003 \times 0,9968}{0,003 \times 0,9968 + (1 - 0,003) \times (1 - 0,9991)} \approx 0,77.$$

c. Sur l'ensemble de la population sud-africaine en 2017, la proportion de personnes infectées par le VIH est $p = \frac{7\,000\,000}{35\,000\,000} = 0,2$ donc la valeur prédictive du TROD salivaire est

$$\frac{0,2 \times 0,9803}{0,2 \times 0,9803 + (1 - 0,2) \times 0,9974} \approx 0,99$$

et la valeur prédictive du TROD sanguin est

$$\frac{0,2 \times 0,9968}{0,2 \times 0,9968 + (1 - 0,2) \times 0,9991} \approx 0,996.$$

d. Les résultats précédents montrent l'influence de la prévalence sur la valeur prédictive positive. On voit bien à travers les trois exemples que plus la prévalence est grande, plus la valeur prédictive l'est aussi.

Exercice 2.

1. D'après l'énoncé, la sensibilité d'une mammographie est $Se = \frac{82}{100} = 0,82$ et la spécificité d'une mammographie est $Sp = 1 - \frac{9}{100} = 0,91$.

2.

	Anomalie détectée	Pas d'anomalie détectée	Total
Personnes malades	82	18	100
Personnes non malades	891	9 009	9 900
Total	973	9 027	10 000

3. La probabilité qu'une femme ayant une anomalie détecté soit atteinte d'un cancer du sein est $\frac{82}{973} \approx 0,084$.

4. La valeur calculée à la question précédente correspond exactement à la valeur prédictive positive. On peut le vérifier avec la formule du cours :

$$\frac{0,01 \times 0,82}{0,01 \times 0,82 + (1 - 0,10)(1 - 0,91)} \approx 0,084.$$

La valeur prédictive négative est la probabilité qu'une personne chez qui aucun anomalie n'est détectée ne soit pas malade. Celle-ci est donc égale à $\frac{9009}{9027} \approx 0,998$.

Exercice 3. Si $Se + Sp = 1$ alors $Sp = 1 - Se$ donc la valeur prédictive du test est

$$\frac{p \times Se}{p \times Se + (1 - p)Se} = \frac{p \times Se}{p \times Se + Se - p \times Se} = \frac{p \times Se}{Se} = p.$$

Ainsi, la valeur prédictive du test est égal à la prévalence donc le test n'a pas d'intérêt car la probabilité qu'un personne prise au hasard soit malade est la même qu'elle ait un test positif ou qu'elle n'est pas fait de test.

Exercice 4.

1. a. Notons D : « la personne est atteinte de la Dengue », Ch : « la personne est atteinte du Chikungunya », Z : « la personne est infectée par le virus Zika », F : « la personne a de la fièvre », Cb : « la personne a des courbatures » et O : « la personne a des douleurs oculaires ». On peut représenter la situation par l'arbre pondéré suivant :
Par indépendance, la probabilité qu'une personne présente de la fièvre et des douleurs oculaires mais pas de courbatures est

$$P(F \cap O \cap \overline{Cb}) = P(F)P(\overline{Cb})P(O) = P(F)(1 - P(Cb))P(O).$$

Or, comme D, Ch et Z forment une partition de l'univers, d'après la formule des probabilités totales,

$$P(F) = P(D)P_D(F) + P(Ch)P_{Ch}(F) + P(Z)P_Z(F) = \frac{1}{3} \times 0,95 + \frac{1}{3} \times 0,75 + \frac{1}{3} \times 0,75 = \frac{49}{60}.$$

On obtient, de même,

$$P(Cb) = \frac{1}{3}(0,75 + 0,95 + 0,5) = \frac{11}{15}$$

et

$$P(O) = \frac{1}{3}(0,5 + 0,25 + 0,5) = \frac{5}{12}.$$

donc

$$P(F \cap \overline{Cb} \cap O) = \frac{49}{60} \times \left(1 - \frac{11}{15}\right) \times \frac{5}{12} = \frac{49}{540}.$$

Toujours par indépendance, si une personne a le Dengue, la probabilité qu'elle présente de la fièvre, des douleurs oculaires mais pas de courbatures est

$$P_D(F \cap \overline{Cb} \cap O) = P_D(F)P_D(1 - P_D(Cb))P_D(O) = 0,95 \times (1 - 0,75) \times 0,5 = \frac{19}{160}.$$

De même, la probabilité qu'elle présente de la fièvre, des douleurs oculaires mais pas de courbatures sachant qu'elle est atteinte du Chikungunya est

$$P_{Ch}(F \cap \overline{Cb} \cap O) = 0,75 \times (1 - 0,95) \times 0,25 = \frac{3}{320}$$

et la probabilité qu'elle présente de la fièvre, des douleurs oculaires mais pas de courbatures sachant qu'elle est infecté par le virus Zika

$$P_Z(F \cap \overline{Cb} \cap O) = 0,75 \times (1 - 0,5) \times 0,5 = \frac{3}{16}.$$

Ainsi, si la personne présente à la fois de la fièvre, pas de courbatures et des douleurs oculaires, la probabilité qu'elle soit atteinte de la Dengue est

$$P_{F \cap \overline{Cb} \cap O}(D) = \frac{P(D) \times P_D(F \cap \overline{Cb} \cap O)}{P(F \cap \overline{Cb} \cap O)} = \frac{\frac{1}{3} \times \frac{19}{160}}{\frac{49}{540}} = \frac{171}{392}.$$

De même, si la personne présente à la fois de la fièvre, pas de courbatures et des douleurs oculaires, la probabilité qu'elle soit atteinte du Chikungunya est

$$P_{F \cap \overline{Cb} \cap O}(Ch) = \frac{\frac{1}{3} \times \frac{3}{320}}{\frac{49}{540}} = \frac{27}{784}$$

et la probabilité qu'elle soit infecté par le virus Zika est

$$P_{F \cap \overline{Cb} \cap O}(Z) = \frac{\frac{1}{3} \times \frac{3}{16}}{\frac{49}{540}} = \frac{135}{196}$$

b. Dans ce cas, le plus probable est que la personne soit infecté par le virus Zika.

2. a. En reprenant la démarche précédente,

$$P(F) = \frac{10}{100} \times 0,95 + \frac{10}{100} \times 0,75 + \frac{80}{100} \times 0,75 = \frac{77}{100},$$

$$P(Cb) = \frac{10}{100} \times 0,75 + \frac{10}{100} \times 0,95 + \frac{80}{100} \times 0,5 = \frac{57}{100}$$

et

$$P(O) = \frac{10}{100} \times 0,5 + \frac{10}{100} \times 0,25 + \frac{80}{100} \times 0,5 = \frac{19}{40}.$$

donc

$$P(F \cap O \cap \overline{Cb}) = \frac{77}{100} \times \left(1 - \frac{57}{100}\right) \times \frac{19}{40} = \frac{62\,909}{400\,000}.$$

Les probabilités conditionnelles $F \cap \overline{Cb} \cap O$ sachant D, Ch ou Z restent inchangées donc si la personne présente à la fois de la fièvre, pas de courbatures et des douleurs oculaires, la probabilité qu'elle soit atteinte de la Dengue est

$$P_{F \cap \overline{Cb} \cap O}(D) = \frac{\frac{10}{100} \times \frac{19}{160}}{\frac{62\,909}{400\,000}} = \frac{250}{3\,311},$$

la probabilité qu'elle soit atteinte du Chikungunya est

$$P_{F \cap \overline{Cb} \cap O}(\text{Ch}) = \frac{\frac{10}{100} \times \frac{3}{320}}{\frac{62\,909}{400\,000}} = \frac{375}{62\,909}$$

et la probabilité qu'elle soit infecté par le virus Zica est

$$P_{F \cap \overline{Cb} \cap O}(Z) = \frac{\frac{80}{100} \times \frac{3}{16}}{\frac{62\,909}{400\,000}} = \frac{60\,000}{62\,909}$$

b. Dans ce cas encore, le plus probable est que la personne soit infecté par le virus Zica.

Exercice 5. Notons p la prévalence de la maladie, M l'évènement « la personne est malade » et T l'évènement « la personne a un test positif ». Alors, par définition, $p = P(M)$, la sensibilité du test est $\text{Se} = P_M(T)$ et la spécificité du test est $\text{Sp} = P_{\overline{M}}(\overline{T})$.

Comme les évènements M et \overline{M} forment une partition de l'univers, par la formule des probabilités totales,

$$P(T) = P(M)P_M(T) + P(\overline{M})P_{\overline{M}}(T) = p \times \text{Se} + (1 - p) \times (1 - \text{Sp}).$$

Par suite, la valeur prédictive positive du test est donc

$$V = P_T(M) = \frac{P(T \cap M)}{P(T)} = \frac{P(M)P_M(T)}{P(T)} = \frac{p \times \text{Se}}{p \times \text{Se} + (1 - p) \times (1 - \text{Sp})}.$$