

Intelligence Artificielle et Machine Learning

Fabien Lauer

fabien.lauer@loria.fr

5A II Polytech

- ▶ 36h (cours/TD/TP)
- ▶ Evaluation : 2 notes

Contenu

- ▶ Notion d'agent intelligent
- ▶ Algorithmes de recherche pour la résolution de problèmes
- ▶ Jeux
- ▶ Modélisation de l'incertain (probabilités, réseaux bayesiens...),
- ▶ Apprentissage supervisé (discrimination, régression)

Bibliographie

- ▶ *Artificial Intelligence: a modern approach* de S. Russel, P. Norvig
- ▶ *The Elements of Statistical Learning* de T. Hastie, R. Tibshirani, J. Friedman

Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

Classification multi-classe

Méthodes non linéaires

Méthodes à noyaux

Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

Classification multi-classe

Méthodes non linéaires

Méthodes à noyaux

Qu'est-ce que l'Intelligence Artificielle (IA) ?

- ▶ IA = étude, analyse et conception de systèmes intelligents

Qu'est-ce que l'Intelligence Artificielle (IA) ?

- ▶ IA = étude, analyse et conception de systèmes intelligents
- ▶ IA = liste de problèmes qu'un "bête" ordinateur ne peut pas résoudre
casse-têtes ; gagner contre le champion du monde aux dames, aux échecs, au go ; garer une voiture ; reconnaître l'écriture, la parole, un chien dans une photo ; comprendre une phrase ; attraper des objets ; écriture de la musique comme Bach...

Qu'est-ce que l'Intelligence Artificielle (IA) ?

- ▶ IA = étude, analyse et conception de systèmes intelligents
- ▶ IA = liste de problèmes qu'un "bête" ordinateur ne peut pas résoudre
casse-têtes ; gagner contre le champion du monde aux dames, aux échecs, au go ; garer une voiture ; reconnaître l'écriture, la parole, un chien dans une photo ; comprendre une phrase ; attraper des objets ; écriture de la musique comme Bach...
- ▶ IA = différentes définitions construites selon le but et qui se distinguent selon deux axes :

Un système intelligent est un système qui est capable...

| | |
|---------------------------|------------------------|
| de penser comme un humain | de raisonner |
| d'agir comme un humain | d'agir raisonnablement |

Agir comme un humain

Test(s) de Turing

- ▶ Est-ce que la machine peut avoir une conversion naturelle avec un humain ?
- ▶ Est-ce qu'un humain peut faire la différence entre une conversation tenue par un humain ou par une machine ?
- ▶ A l'écrit, en face à face, en autorisant le passage d'objets ?

- ▶ Pas si intéressant que ça... tant qu'on ne souhaite pas interagir directement avec (ou remplacer) un humain
- ▶ Mais inclut beaucoup de problèmes d'IA
 - ▶ traitement du langage naturel
 - ▶ représentation des connaissances
 - ▶ raisonnement automatique
 - ▶ apprentissage automatique
 - ▶ vision (pour le face à face)
 - ▶ robotique

- ▶ Cleverbot (en 2011, "59% humain" contre 63% pour les vrais humains)

Penser comme un humain

Sciences cognitives

- ▶ Comprendre le fonctionnement du cerveau
- ▶ En le simulant avec une machine
- ▶ Si les résultats de la simulation coïncident avec les observations sur les humains, alors le modèle implanté nous apprend peut-être quelque chose sur le cerveau

- ▶ La machine n'est qu'un outil de simulation, pas une fin en soi
- ▶ On ne cherche pas à résoudre des problèmes ou à "remplacer l'humain"

- ▶ Human Brain Project (projet européen d'envergure, 2014–2024)

Logique

- ▶ Implémenter les règles de la pensée rationnelle dans les machines
- ▶ Déductions automatiques à partir d'axiomes / de faits
- ▶ Limites / difficultés
 - ▶ Formulation des connaissances informelles en langage formel
 - ▶ Les connaissances sont souvent approximatives plutôt que des faits certains
 - ▶ La complexité des problèmes devient vite énorme (explosion du nombre de raisonnements possibles à partir de quelques dizaines de faits)

Agir raisonnablement

Théorie de la décision, statistique, apprentissage

- ▶ Agir le mieux possible pour atteindre un objectif à partir d'une certaine connaissance/de certaines croyances sur le monde
- ▶ Décider qu'une action est meilleure dans certaines circonstances
- ▶ Avantages
 - ▶ Cette définition est claire et précise
 - ▶ Les problèmes peuvent se formuler mathématiquement
 - ▶ et être analysés en s'appuyant sur beaucoup d'outils pour construire une théorie solide
 - ▶ Nombreuses applications pratiques
- ▶ Inconvénients
 - ▶ Ne raisonne pas forcément comme un humain
 - ▶ Ne permet pas de comprendre le cerveau

Agir raisonnablement

Théorie de la décision, statistique, apprentissage

- ▶ Agir le mieux possible pour atteindre un objectif à partir d'une certaine connaissance/de certaines croyances sur le monde
- ▶ Décider qu'une action est meilleure dans certaines circonstances
- ▶ Avantages
 - ▶ Cette définition est claire et précise
 - ▶ Les problèmes peuvent se formuler mathématiquement
 - ▶ et être analysés en s'appuyant sur beaucoup d'outils pour construire une théorie solide
 - ▶ Nombreuses applications pratiques
- ▶ Inconvénients
 - ▶ Ne raisonne pas forcément comme un humain
 - ▶ Ne permet pas de comprendre le cerveau

Approche suivie dans ce cours

Qu'est-ce que l'Intelligence Artificielle (IA) ?

- ▶ IA = étude, analyse et conception de systèmes intelligents
- ▶ IA = liste de problèmes qu'un "bête" ordinateur ne peut pas résoudre
casse-têtes ; gagner contre le champion du monde aux dames, aux échecs, au go ; garer une voiture ; reconnaître l'écriture, la parole, un chien dans une photo ; comprendre une phrase ; attraper des objets ; écriture de la musique comme Bach...
- ▶ IA = différentes définitions construites selon le but et qui se distinguent selon deux axes :

Un système intelligent est un système qui est capable...

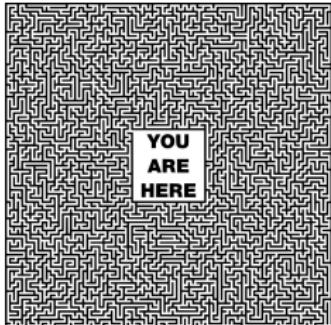
| | |
|---------------------------|------------------------|
| de penser comme un humain | de raisonner |
| d'agir comme un humain | d'agir raisonnablement |

- ▶ Aujourd'hui : IA = une liste de mots clés :
search, problem-solving, machine learning, data mining, artificial neural networks, deep learning, pattern recognition, statistics, case-based reasoning, logic, ... big data, image recognition/retrieval, speech synthesis/recognition, natural language processing, automated translation, autonomous cars, robotics, games, planning, path-finding.....

Qu'est-ce que l'Intelligence Artificielle (IA) ?

IA des années 50 aux années 80

- ▶ Résoudre des problèmes complexes pour les humains avec des algorithmes simples



comment sortir ?

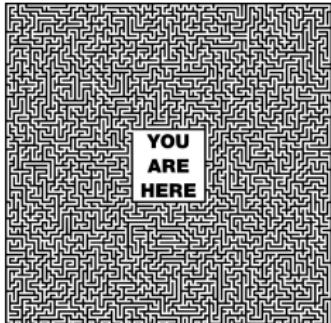


que faire ensuite ?

Qu'est-ce que l'Intelligence Artificielle (IA) ?

IA des années 50 aux années 80

- ▶ Résoudre des problèmes complexes pour les humains avec des algorithmes simples



comment sortir ?



que faire ensuite ?

A partir des années 80 : machine learning

- ▶ Calculer des choses évidentes avec des algorithmes complexes

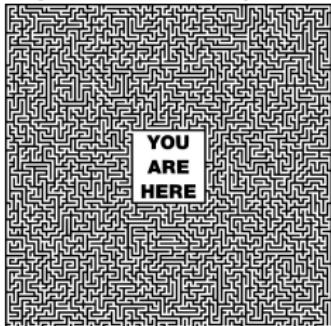


Est-ce un chat ou un chien ?

Qu'est-ce que l'Intelligence Artificielle (IA) ?

IA des années 50 aux années 80

- ▶ Résoudre des problèmes complexes pour les humains avec des algorithmes simples



comment sortir ?



que faire ensuite ?

A partir des années 80 : machine learning

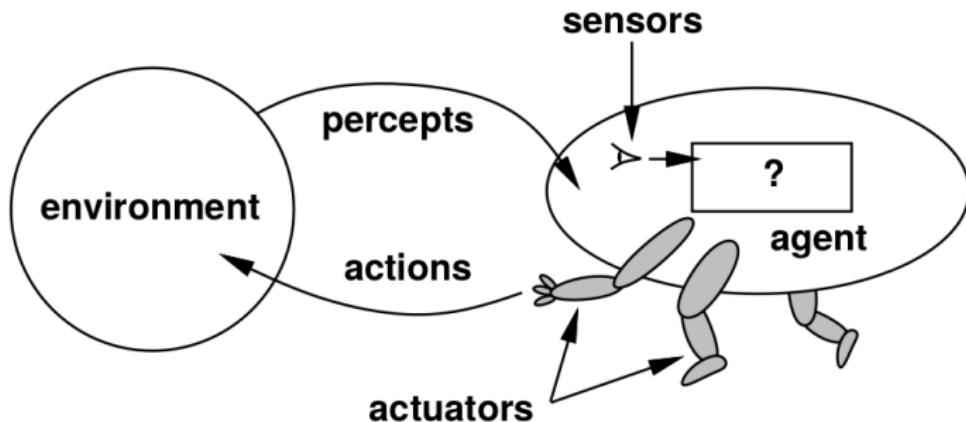
- ▶ Calculer des choses évidentes avec des algorithmes complexes



Est-ce un chat ou un chien ?

+ d'autres choses un peu moins évidentes...

Agents intelligents



Un agent...

- ▶ vie/évolue dans un **environnement**
- ▶ **perçoit** son environnement grâce à ses **capteurs**
- ▶ **agit** sur cet environnement grâce à des actionneurs/effecteurs

Agent rationnel = agent qui exécute la meilleure action possible étant donné sa perception du monde

Mesure de performance

Un agent (ou une séquence d'actions) est "meilleur" qu'un autre = sa mesure de performance est plus grande

Exemples de mesures de performances

La performance d'un aspirateur intelligent se mesure en fonction de

- ▶ la quantité de poussière aspirée ?

Mesure de performance

Un agent (ou une séquence d'actions) est "meilleur" qu'un autre = sa mesure de performance est plus grande

Exemples de mesures de performances

La performance d'un aspirateur intelligent se mesure en fonction de

- ▶ la quantité de poussière aspirée ?
- ▶ et l'électricité consommée ?

Mesure de performance

Un agent (ou une séquence d'actions) est "meilleur" qu'un autre = sa mesure de performance est plus grande

Exemples de mesures de performances

La performance d'un aspirateur intelligent se mesure en fonction de

- ▶ la quantité de poussière aspirée ?
- ▶ et l'électricité consommée ?
- ▶ et le niveau de bruit ?

Mesure de performance

Un agent (ou une séquence d'actions) est "meilleur" qu'un autre = sa mesure de performance est plus grande

Exemples de mesures de performances

La performance d'un aspirateur intelligent se mesure en fonction de

- ▶ la quantité de poussière aspirée ?
- ▶ et l'électricité consommée ?
- ▶ et le niveau de bruit ?

- ▶ sur une heure ? sur la journée ? sur l'année ?

Mesure de performance

Un agent (ou une séquence d'actions) est "meilleur" qu'un autre = sa mesure de performance est plus grande

Exemples de mesures de performances

La performance d'un aspirateur intelligent se mesure en fonction de

- ▶ la quantité de poussière aspirée ?
- ▶ et l'électricité consommée ?
- ▶ et le niveau de bruit ?

- ▶ sur une heure ? sur la journée ? sur l'année ?

Construire une mesure de performance peut être complexe, mais constitue une des étapes les plus importantes (car cela détermine directement ce qu'un agent rationnel fera)

Limitations dues aux perceptions

Agent rationnel vs omniscient

- ▶ Un agent rationnel peut décider de traverser la rue et mourir en se prenant un frigo tombé d'un avion sur la tête
- ▶ Car ses perceptions ne lui permettent pas de savoir que le frigo arrive
- ▶ Son action était rationnelle = idéale étant donné ses connaissances
- ▶ MAIS traverser la rue sans regarder à gauche et à droite n'est pas rationnel
Ne pas oublier qu'acquérir de l'information peut être une action

Agent rationnel idéal

- ▶ Pour toute séquence de perception, prend l'action qui **maximise la performance** attendue sur la base des perceptions et de ses connaissances
- ▶ Pour pouvoir s'adapter, un agent doit être **autonome** : ses actions ne dépendent pas que de son constructeur mais aussi de son expérience

Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

Agents qui résolvent des problèmes

Formulation d'un problème et d'un objectif

- ▶ Mesure de performance : inclue beaucoup de critères

Je suis perdu en Roumanie et...

- ▶ je dois prendre l'avion à Bucarest dans 2 jours
- ▶ je souhaiterais visiter des sites touristiques
- ▶ je souhaiterais voir un peu la campagne
- ▶ je n'aime pas trop les bosses sur la route
- ▶ j'aime bien écouter la radio en roulant
- ▶ je n'ai pas trop d'argent pour l'essence...

Agents qui résolvent des problèmes

Formulation d'un problème et d'un objectif

- ▶ Mesure de performance : inclue beaucoup de critères
- ▶ Dans certaines situations, il faut se concentrer sur 1 objectif précis
- ▶ **But** = ensemble des états dans lesquels l'objectif est atteint
- ▶ Solution = suite d'actions menant au but

Je suis perdu en Roumanie et...

- ▶ je dois **prendre l'avion à Bucarest** dans 2 jours
- ▶ je souhaiterais visiter des sites touristiques
- ▶ je souhaiterais voir un peu la campagne
- ▶ je n'aime pas trop les bosses sur la route
- ▶ j'aime bien écouter la radio en roulant
- ▶ je n'ai pas trop d'argent pour l'essence...

Agents qui résolvent des problèmes

Formulation d'un problème et d'un objectif

- ▶ Mesure de performance : inclue beaucoup de critères
- ▶ Dans certaines situations, il faut se concentrer sur 1 objectif précis
- ▶ **But** = ensemble des états dans lesquels l'objectif est atteint
- ▶ Solution = suite d'actions menant au but
- ▶ **Coût** d'une solution : certains problèmes ont plusieurs solutions, certaines ont un coût plus élevé, on cherchera donc la solution avec le coût le plus faible

Un agent doit trouver une séquence d'actions qui le mène au but = recherche puis l'exécute = exécution

Je suis perdu en Roumanie et...

- ▶ je dois **prendre l'avion à Bucarest dans 2 jours (~ le plus vite possible)**
- ▶ je souhaiterais visiter des sites touristiques
- ▶ je souhaiterais voir un peu la campagne
- ▶ je n'aime pas trop les bosses sur la route
- ▶ j'aime bien écouter la radio en roulant
- ▶ je n'ai pas trop d'argent pour l'essence...

Agents qui résolvent des problèmes

Différents cas de figure

- ▶ Cas 1 : monde accessible et effets des actions connus
- ▶ Cas 2 : monde partiellement accessible et effets des actions connus
- ▶ Cas 3 : monde accessible et effets des actions connus partiellement
- ▶ Cas 4 : monde partiellement accessible et effets des actions partiellement connus

Exemple de l'aspirateur...

Agents qui résolvent des problèmes

Différents cas de figure

- ▶ Cas 1 : monde accessible et effets des actions connus
- ▶ Cas 2 : monde partiellement accessible et effets des actions connus
- ▶ Cas 3 : monde accessible et effets des actions connus partiellement
- ▶ Cas 4 : monde partiellement accessible et effets des actions partiellement connus

Solutions

- ▶ Cas 1 : Séquence d'actions = séquence d'états \Rightarrow trouver celle qui conduit au but
- ▶ Cas 2 : Idem mais avec une séquence d'ensembles d'états possibles
- ▶ Cas 3 : Séquence d'actions conditionnées aux perceptions
- ▶ Cas 4 : Pas toujours de solution figée qui garantisse d'arriver au but... il faut itérer entre la recherche d'une solution (partielle) et l'exécution d'actions

Exemple de l'aspirateur...

Algorithmes de recherche et résolution de problèmes

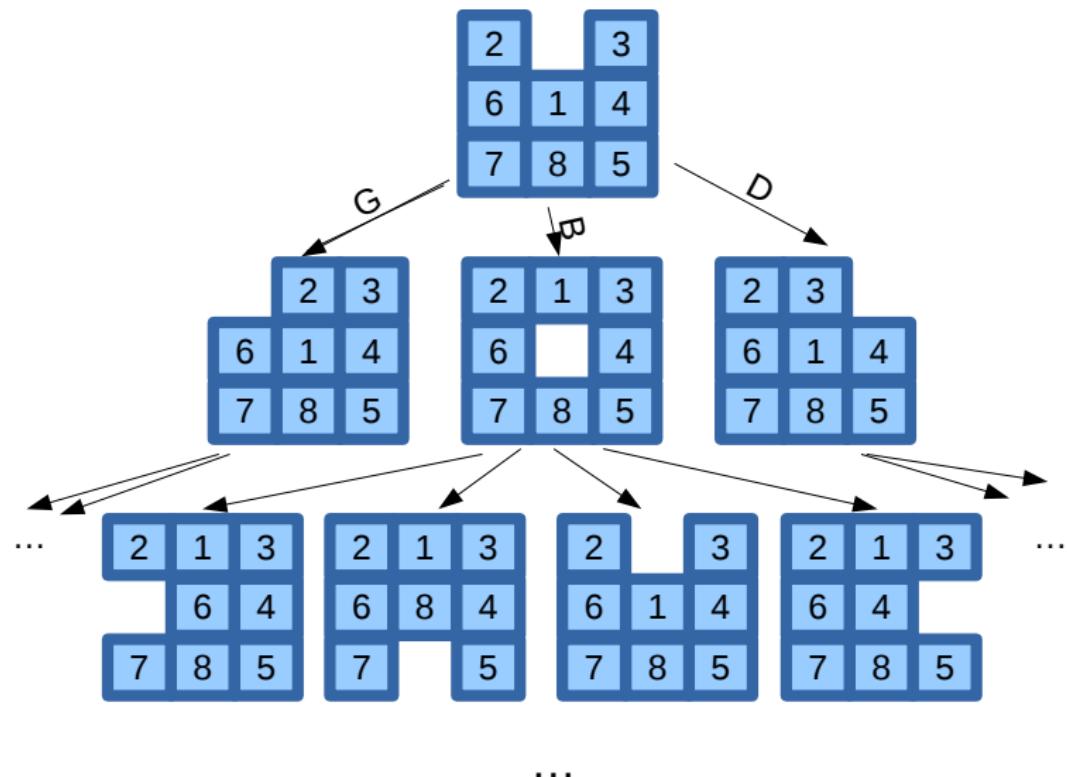
Arbre de recherche

- ▶ Racine = état initial
- ▶ 1 nœud pour chaque état du problème
- ▶ 1 branche pour chaque action possible, pondérée par le coût de l'action
- ▶ On descend dans l'arbre en suivant une séquence d'actions jusqu'à un nœud correspondant au but (ensemble d'états satisfaisant l'objectif)
- ▶ Le coût d'une solution est la somme des pondérations des branches traversées

Exemple : taquin

- ▶ Etat : position des pièces et du trou
- ▶ Actions possibles : déplacer le trou vers le haut, le bas, la gauche ou la droite
(en fonction de l'état, certaines actions sont impossibles)
- ▶ Coût d'une action = 1 pour toutes les actions
on cherche à minimiser le nombre de coups pour arriver au but

Arbre de recherche



Algorithmes de recherche et résolution de problèmes

Algorithmes de recherche "aveugles"

- ▶ **Largeur d'abord** : garantie de trouver une solution (avec le moins d'actions possibles) mais demande beaucoup de mémoire
- ▶ **Coût uniforme** : idem largeur d'abord, mais trouve la solution de coût minimal (mais peut-être plus d'actions) en creusant les nœuds qui minimisent le coût $g(\text{noeud})$ entre la racine et le nœud
- ▶ **Profondeur d'abord** : plus efficace en terme de mémoire, peut trouver une solution rapidement mais peut aussi perdre beaucoup de temps dans des branches inintéressantes
- ▶ Mélange : augmenter progressivement la profondeur maximale d'une recherche en profondeur d'abord

Algorithmes de recherche et résolution de problèmes

Algorithmes de recherche "aveugles"

- ▶ **Largeur d'abord** : garantie de trouver une solution (avec le moins d'actions possibles) mais demande beaucoup de mémoire
- ▶ **Coût uniforme** : idem largeur d'abord, mais trouve la solution de coût minimal (mais peut-être plus d'actions) en creusant les nœuds qui minimisent le coût $g(\text{noeud})$ entre la racine et le nœud
- ▶ **Profondeur d'abord** : plus efficace en terme de mémoire, peut trouver une solution rapidement mais peut aussi perdre beaucoup de temps dans des branches inintéressantes
- ▶ Mélange : augmenter progressivement la profondeur maximale d'une recherche en profondeur d'abord

Algorithmes de recherche "informés" : le meilleur d'abord

Avec une heuristique $h(\text{noeud})$ qui estime le coût restant pour le meilleur chemin entre *noeud* et *but*

- ▶ **Recherche gloutonne** : creuser le nœud qui minimise $h(\text{noeud})$ pour se rapprocher le plus possible du but ; mais sans garantie sur le coût de la solution
- ▶ **A*** : combiner coût uniforme et recherche gloutonne : creuser le nœud qui minimise $f(\text{noeud}) = g(\text{noeud}) + h(\text{noeud})$ = estimation du coût total d'une solution passant par *noeud*

A^* : heuristiques admissibles

- ▶ Une heuristique est dite **admissible** si $h(\text{noeud}) \leq$ coût du meilleur chemin de *noeud* à *but* pour tout *noeud*
- ▶ Si l'heuristique est admissible alors A^* est garanti de trouver une **solution optimale**

A^* : heuristiques admissibles

- ▶ Une heuristique est dite **admissible** si $h(\text{noeud}) \leq$ coût du meilleur chemin de *noeud* à *but* pour tout *noeud*
- ▶ Si l'heuristique est admissible alors A^* est garanti de trouver une **solution optimale**

Problème du taquin

- ▶ $g(n) =$ nombre de coups joués
- ▶ $h_1(n) =$ nombre de pièces mal placées
- ▶ $h_2(n) =$ somme des distances des pièces à leur position finale
- ▶ h_1 et h_2 admissibles car nb coups qu'il reste à jouer $\geq h_1(n)$ et $\geq h_2(n)$
- ▶ Quelle est la meilleure heuristique ?

A^* : heuristiques admissibles

- ▶ Une heuristique est dite **admissible** si $h(\text{noeud}) \leq$ coût du meilleur chemin de *noeud* à *but* pour tout *noeud*
- ▶ Si l'heuristique est admissible alors A^* est garanti de trouver une **solution optimale**

Problème du taquin

- ▶ $g(n) =$ nombre de coups joués
- ▶ $h_1(n) =$ nombre de pièces mal placées
- ▶ $h_2(n) =$ somme des distances des pièces à leur position finale
- ▶ h_1 et h_2 admissibles car nb coups qu'il reste à jouer $\geq h_1(n)$ et $\geq h_2(n)$
- ▶ Quelle est la meilleure heuristique ?
celle qui estime le mieux le coût restant
- ▶ Puisque coût restant \geq heuristique, celle qui estime le mieux le coût restant est celle qui est la plus grande *pour tous les nœuds*
- ▶ $h_2(n) \geq h_1(n)$ pour tout n donc h_2 est meilleure !

Algorithmes de recherche en pratique

Algorithme générique

Squelette identique pour tous les algorithmes de recherche :

1. Prendre un nœud dans la file d'attente
2. Vérifier si le but est atteint à ce nœud (si oui, retourner la solution)
3. Développer le nœud et ajouter les fils à la file d'attente

Un algorithme de recherche = algorithme générique + fonction pour ordonner la file d'attente

- ▶ Largeur d'abord : ajouter les nouveaux nœuds à la fin de la file
- ▶ Coût uniforme : trier la file par ordre croissant de $g(\text{noeud})$
- ▶ Profondeur d'abord : ajouter les nouveaux nœuds au début de la file
- ▶ Glouton : trier la file par ordre croissant de $h(\text{noeud})$
- ▶ A^* : trier la file par ordre croissant de $f(\text{noeud}) = g(\text{noeud}) + h(\text{noeud})$

Eviter les états répétés

Problème

- ▶ Plusieurs chemins mènent au même état, donc dans un arbre le même état est évalué plusieurs fois...
- ▶ En profondeur d'abord, ou glouton, encore plus grave : problème de profondeur infinie
- ▶ Pour les autres, limiter les répétitions accélère aussi l'algorithme

Solutions

3 niveaux de plus en plus efficaces mais de plus en plus coûteux

1. Ne pas revenir en arrière
2. Ne pas générer un état présent dans le chemin qui y mène (pour éviter les cycles)
3. Ne pas générer d'état présent dans l'arbre

En règle générale en pratique : solutions 1 & 2

Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

Qu'est-ce qu'un jeu ?

Définitions / limitation du cadre d'étude

- ▶ Jeux déterministes (dans un premier temps) : pas de hasard
- ▶ Jeux à information parfaite : tous les éléments sont connus de tous les joueurs
- ▶ Jeux à 2 joueurs (un contre l'autre)

Remarque : les jeux à plus de 2 joueurs sont forcément à information imparfaite (alliances...)

Qu'est-ce qu'un jeu ?

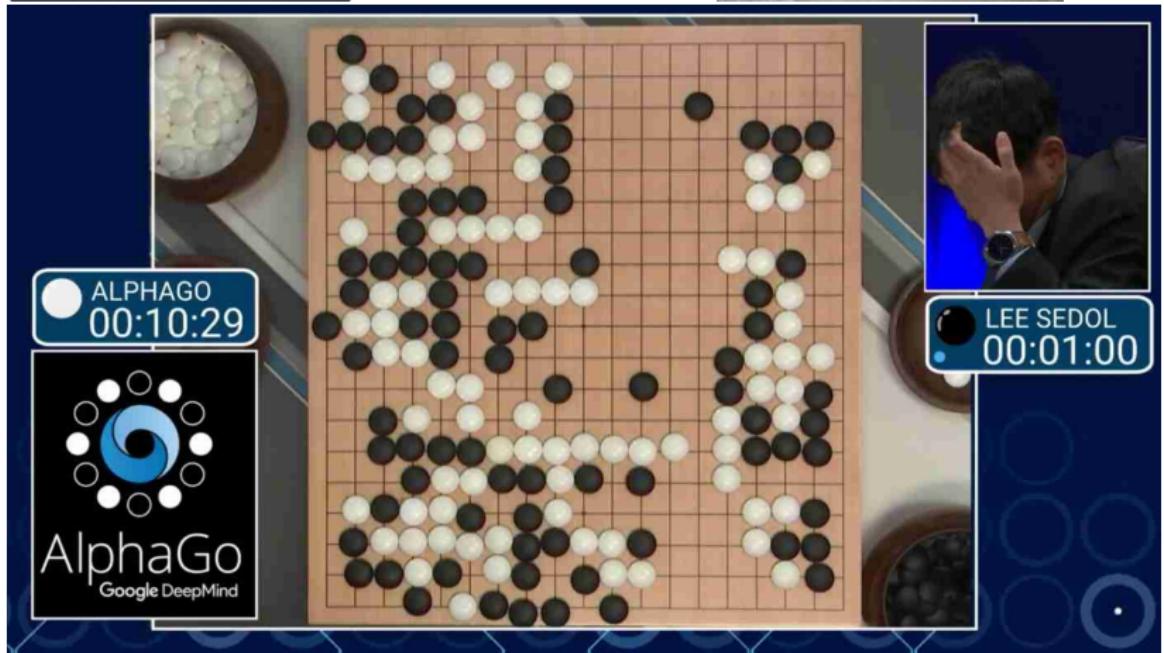
Définitions / limitation du cadre d'étude

- ▶ Jeux déterministes (dans un premier temps) : pas de hasard
- ▶ Jeux à information parfaite : tous les éléments sont connus de tous les joueurs
- ▶ Jeux à 2 joueurs (un contre l'autre)

Remarque : les jeux à plus de 2 joueurs sont forcément à information imparfaite (alliances...)

Exemples

- ▶ Dames : Chinook (champion du monde en 1994) (alpha-beta + répertoire de fins de partie)
- ▶ Othello/reversi
- ▶ Echecs : Kasparov vs IBM DeepBlue (1997) (alpha-beta + répertoires d'ouvertures et de fins + machine parallèle dédiée)
- ▶ Jeopardy : IBM Watson (2011) (rien à voir...)
- ▶ Go : AlphaGO (champion européen en 2013, champion du monde en 2015) (MCTS + apprentissage)



Jeux à deux joueurs (Max vs Mini)

Formulation d'un jeu

- ▶ **Etat initial** : position des pièces, à qui de jouer...
- ▶ Ensemble d'**opérateurs** : définit tous les coups possibles pour un joueur (fonction de transition d'un état à l'autre)
- ▶ **Test de terminaison** : est-ce que la partie est finie ?
- ▶ **Fonction de récompense** : retourne une valeur numérique à la fin de la partie en fonction de l'état du jeu (par ex. $\{+1, 0, -1\}$ ou un score)

Jeux à deux joueurs (Max vs Mini)

Formulation d'un jeu

- ▶ **Etat initial** : position des pièces, à qui de jouer...
- ▶ Ensemble d'**opérateurs** : définit tous les coups possibles pour un joueur (fonction de transition d'un état à l'autre)
- ▶ **Test de terminaison** : est-ce que la partie est finie ?
- ▶ **Fonction de récompense** : retourne une valeur numérique à la fin de la partie en fonction de l'état du jeu (par ex. $\{+1, 0, -1\}$ ou un score)

Arbre de jeu

- ▶ Chaque nœud correspond à un état du jeu
- ▶ Chaque branche correspond à un coup possible
- ▶ Chaque niveau/profondeur est associé au joueur qui doit jouer

Jeux à deux joueurs (Max vs Mini)

Formulation d'un jeu

- ▶ **Etat initial** : position des pièces, à qui de jouer...
- ▶ Ensemble d'**opérateurs** : définit tous les coups possibles pour un joueur (fonction de transition d'un état à l'autre)
- ▶ **Test de terminaison** : est-ce que la partie est finie ?
- ▶ **Fonction de récompense** : retourne une valeur numérique à la fin de la partie en fonction de l'état du jeu (par ex. $\{+1, 0, -1\}$ ou un score)

Arbre de jeu

- ▶ Chaque nœud correspond à un état du jeu
- ▶ Chaque branche correspond à un coup possible
- ▶ Chaque niveau/profondeur est associé au joueur qui doit jouer

Différent d'un problème de recherche car l'adversaire agit sur le résultat : il ne suffit pas de prendre un chemin qui mène à la meilleure récompense, mais un chemin qui mène à la bonne récompense quelles que soient les actions de l'adversaire

Minimax

Action optimale dans un jeu à deux joueurs

- ▶ Sous l'hypothèse que l'autre joueur joue de manière optimale
- ▶ Avec une information complète sur l'état du jeu
- ▶ En énumérant toutes les séquences de jeu possibles (en développant tout l'arbre de jeu)
- ▶ En associant une valeur à chaque nœud

Minimax

Action optimale dans un jeu à deux joueurs

- ▶ Sous l'hypothèse que l'autre joueur joue de manière optimale
- ▶ Avec une information complète sur l'état du jeu
- ▶ En énumérant toutes les séquences de jeu possibles (en développant tout l'arbre de jeu)
- ▶ En associant une valeur à chaque nœud

Algorithme Minimax (du point de vue de Max)

- ▶ Générer entièrement l'arbre de jeu
- ▶ Evaluer la récompense pour chaque nœud terminal
valeur d'un nœud terminal = récompense
- ▶ Remonter l'arbre en propageant
 - ▶ la **valeur minimale** sur l'ensemble des fils si **c'est à Mini de jouer**
(on admet que l'adversaire joue son meilleur coup)
 - ▶ la **valeur maximale** sur l'ensemble des fils si **c'est à Max de jouer**
(on admet que l'on joue notre meilleur coup)
- ▶ Jouer le coup correspondant au fils de la racine avec la valeur maximale

Minimax

Action optimale dans un jeu à deux joueurs

- ▶ Sous l'hypothèse que l'autre joueur joue de manière optimale
- ▶ Avec une information complète sur l'état du jeu
- ▶ En énumérant toutes les séquences de jeu possibles (en développant tout l'arbre de jeu)
- ▶ En associant une valeur à chaque nœud

Algorithme Minimax (du point de vue de Max)

- ▶ Générer entièrement l'arbre de jeu
- ▶ Evaluer la récompense pour chaque nœud terminal
valeur d'un nœud terminal = récompense
- ▶ Remonter l'arbre en propageant
 - ▶ la **valeur minimale** sur l'ensemble des fils si **c'est à Mini de jouer**
(on admet que l'adversaire joue son meilleur coup)
 - ▶ la **valeur maximale** sur l'ensemble des fils si **c'est à Max de jouer**
(on admet que l'on joue notre meilleur coup)
- ▶ Jouer le coup correspondant au fils de la racine avec la valeur maximale

En pratique : algorithme récursif avec

$$\text{valeurMax}(\text{noeud}) = \max_{f \in \text{fils}(\text{noeud})} \text{valeurMin}(f)$$

$$\text{valeurMin}(\text{noeud}) = \min_{f \in \text{fils}(\text{noeud})} \text{valeurMax}(f)$$

Minimax

Action optimale dans un jeu à deux joueurs

- ▶ Sous l'hypothèse que l'autre joueur joue de manière optimale
- ▶ Avec une information complète sur l'état du jeu
- ▶ En énumérant toutes les séquences de jeu possibles (en développant tout l'arbre de jeu)
- ▶ En associant une valeur à chaque nœud

Algorithme Minimax (du point de vue de Max)

- ▶ Générer entièrement l'arbre de jeu
- ▶ Evaluer la récompense pour chaque nœud terminal
valeur d'un nœud terminal = récompense
- ▶ Remonter l'arbre en propageant
 - ▶ la **valeur minimale** sur l'ensemble des fils si **c'est à Mini de jouer**
(on admet que l'adversaire joue son meilleur coup)
 - ▶ la **valeur maximale** sur l'ensemble des fils si **c'est à Max de jouer**
(on admet que l'on joue notre meilleur coup)
- ▶ Jouer le coup correspondant au fils de la racine avec la valeur maximale

condition d'arrêt :
arriver à une feuille

En pratique : algorithme récursif avec

$$\begin{aligned} valeurMax(noeud) &= \max_{f \in fils(noeud)} valeurMin(f) \quad (\text{ou récompense pour un noeud terminal}) \\ valeurMin(noeud) &= \min_{f \in fils(noeud)} valeurMax(f) \end{aligned}$$

Minimax

Complexité du minimax

- ▶ b : facteur de branchement moyen
 - ~ 10 au reversi, ~ 35 aux échecs, ~ 360 au Go
- ▶ n : profondeur moyenne d'une branche
 - 60 au reversi, ~ 40 aux échecs, ~ 400 au Go
- ▶ Complexité en temps : $O(b^n)$
- ▶ Complexité en espace : $O(bn)$
 - En pratique, avec une approche en profondeur d'abord il est suffisant de mémoriser uniquement la branche en cours d'exploration

Deux solutions possibles

- ▶ Réduire b
- ▶ Réduire n

Alpha–beta (pour réduire b)

Elaguer l'arbre pour gagner du temps

- ▶ L'arbre est parcouru en profondeur d'abord \Rightarrow la valeur des nœuds d'une branche peut être calculée avant de passer à la branche suivante
- ▶ Ne pas explorer les branches qui ne seront jamais jouées car elles ne correspondent pas à des stratégies optimales pour les joueurs

Alpha–beta (pour réduire b)

Elaguer l'arbre pour gagner du temps

- ▶ L'arbre est parcouru en profondeur d'abord \Rightarrow la valeur des nœuds d'une branche peut être calculée avant de passer à la branche suivante
- ▶ Ne pas explorer les branches qui ne seront jamais jouées car elles ne correspondent pas à des stratégies optimales pour les joueurs

Algorithme Alpha-beta : minimax modifié (init : $\alpha = -\infty, \beta = \infty$)

Fonction **valeurMax** (nœud, α, β)

- ▶ $v = -\infty$
- ▶ Pour chaque *fils* de nœud
 - ▶ calculer sa valeur et le max courant $v = \max(v, \text{valeurMin}(\text{fils}, \alpha, \beta))$
 - ▶ Si $v \geq \beta$, retourner v
 - ▶ $\alpha = \max(\alpha, v)$
- ▶ Retourner v

Fonction **valeurMin** (nœud, α, β)

- ▶ $v = +\infty$
- ▶ Pour chaque *fils* de nœud
 - ▶ calculer sa valeur et le min courant $v = \min(v, \text{valeurMax}(\text{fils}, \alpha, \beta))$
 - ▶ Si $v \leq \alpha$, retourner v
 - ▶ $\beta = \min(\beta, v)$
- ▶ Retourner v

Fonctions d'évaluation et limitation de profondeur

Réduire la profondeur n

- ▶ Limiter la profondeur de recherche à n (par exemple $n = 5$ ou $n = 10$)
- ▶ Estimer la qualité des nœuds à la profondeur n par une fonction d'évaluation plutôt que de calculer la valeur exacte du nœud à partir de la récompense finale

Fonction d'évaluation

- ▶ Doit correspondre à la récompense si appliquée à une position terminale
- ▶ Doit être calculable rapidement
- ▶ Doit donner une bonne image des "chances" de gagner à partir d'une position
- ▶ Pour les jeux "à somme nulle", doit être de la forme **(position symétrique au départ)**
 $\text{evalJoueur}(\text{Max}) - \text{evalJoueur}(\text{Mini})$

Attention : en limitant la profondeur de la recherche, l'algorithme Minimax n'est plus optimal ; il est totalement aveugle à ce qui se passe après n coups

Fonctions d'évaluation et limitation de profondeur

Fonction d'évaluation linéaire

$$\text{evaluation}(\text{etat}) = \sum_{k=1}^p w_k \phi_k(\text{etat})$$

Ex. aux échecs : $\phi_k(\text{etat})$ = nb de pièces du type k , w_k = valeurs des pièces

- ▶ les fonctions de caractéristiques ϕ_k dépendent du jeu...
- ▶ les poids w_k peuvent être déterminés automatiquement par "apprentissage" si un expert évalue un ensemble de positions

Fonctions d'évaluation et limitation de profondeur

Fonction d'évaluation linéaire

$$\text{evaluation}(\text{etat}) = \sum_{k=1}^p w_k \phi_k(\text{etat})$$

Ex. aux échecs : $\phi_k(\text{etat})$ = nb de pièces du type k , w_k = valeurs des pièces

- ▶ les fonctions de caractéristiques ϕ_k dépendent du jeu...
- ▶ les poids w_k peuvent être déterminés automatiquement par "apprentissage" si un expert évalue un ensemble de positions

Jouer contre soi-même pour s'améliorer et apprentissage

- ▶ On peut "apprendre" la fonction d'évaluation à partir d'un historique de parties
- ▶ Pour chaque position on note l'évaluation donnée par la fraction de parties gagnées à partir de cette position
- ▶ Un modèle est entraîné à reproduire ces évaluations (et à généraliser !)
- ▶ Cela nécessite un GRAND nombre de parties pour explorer suffisamment de positions suffisamment de fois
- ▶ En faisant jouer le programme contre lui-même, un grand nombre de parties peuvent être jouées rapidement !
- ▶ Attention : si le programme est mauvais au départ, les parties ne seront pas représentatives des vrais adversaires
- ▶ Solution : commencer par "apprendre" à imiter les joueurs humains

Jeux avec du hasard

- ▶ Un coup ne donne le résultat prévu qu'avec une certaine probabilité
- ▶ Un tirage aléatoire pré-conditionne chaque coup

Jeux avec du hasard

- ▶ Un coup ne donne le résultat prévu qu'avec une certaine probabilité
- ▶ Un tirage aléatoire pré-conditionne chaque coup

Arbre probabiliste : arbre de jeu contenant des "nœuds chance"

- ▶ nœud chance : nœud intermédiaire dont l'action est aléatoire
- ▶ En générale, toutes les positions atteignables à partir d'un état sont atteintes avec une probabilité non nulle
⇒ il faut explorer toutes les branches (pas d'équivalent d' α - β)
- ▶ Une séquence d'actions ne conduit qu'à une certaine probabilité de gagner

Jeux avec du hasard

- ▶ Un coup ne donne le résultat prévu qu'avec une certaine probabilité
- ▶ Un tirage aléatoire pré-conditionne chaque coup

Arbre probabiliste : arbre de jeu contenant des "nœuds chance"

- ▶ nœud chance : nœud intermédiaire dont l'action est aléatoire
- ▶ En générale, toutes les positions atteignables à partir d'un état sont atteintes avec une probabilité non nulle
⇒ il faut explorer toutes les branches (pas d'équivalent d' α - β)
- ▶ Une séquence d'actions ne conduit qu'à une certaine probabilité de gagner

Algorithme **Expectimax**

- ▶ Algorithme Minimax modifié tel que la valeur remontée soit
 - ▶ la valeur minimale sur l'ensemble des fils si c'est à Mini de jouer
 - ▶ la valeur maximale sur l'ensemble des fils si c'est à Max de jouer
 - ▶ la **valeur moyenne** sur l'ensemble des fils si c'est un **nœud chance** :

$$\sum_{i \in fils(n)} P(i) valeurExpectimax(i), \quad P(i) : \text{probabilité d'avoir } i \text{ après } n$$

Algorithmes à base de marches aléatoires

Motivations

- ▶ Pas de fonction d'évaluation disponible (par ex. Go)
- ▶ Réflexion en temps limité (ne pas dépasser + profiter au max du temps)

Algorithme aléatoire uniforme

- ▶ Chaque nœud est évalué selon la probabilité de gagner à partir de là
- ▶ On joue le coup qui a la plus grande probabilité de gagner
- ▶ Les probabilités sont estimées par simulation :
Pour chaque coup possible, générer N parties aléatoires et estimer $P(\text{gagner})$ par la fraction de parties gagnées
- ▶ Une partie aléatoire = choix aléatoire (uniforme) à chaque coup
- ▶ Remarque : les parties sont simulées jusqu'à la fin (nœud terminal)
- ▶ Amélioration : modifier le choix aléatoire en introduisant des connaissances sur les bons (ou mauvais) coups
par ex. toujours choisir un coup gagnant quand cela est possible

Recherche d'arbre Monte Carlo

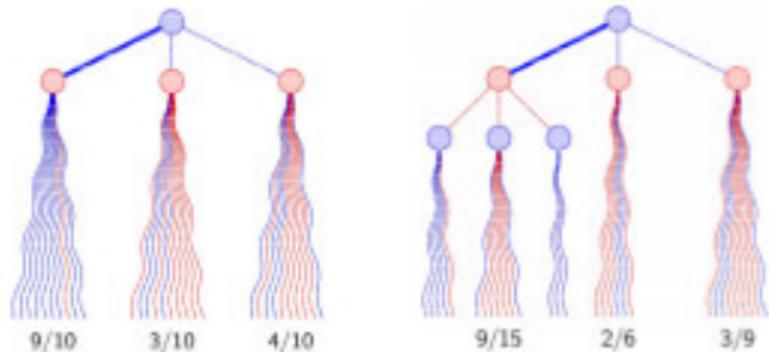
Motivations

- ▶ La recherche aléatoire uniforme peut facilement passer à côté des bons coups et sous-estimer les probabilités de gagner
- ▶ Idée 1 : Explorer les coups qui n'ont pas été suffisamment simulés pour éviter de passer à côté de quelque chose
- ▶ Idée 2 : Exploiter l'information acquise au cours du processus pour améliorer les simulations suivantes et l'estimation des probabilités
- ▶ Chaque idée demande du temps ⇒ **compromis Exploitation / Exploration**

MCTS (*Monte-Carlo Tree Search*)

- ▶ Approche optimiste : on considère que les probabilités sont bien estimées avant la fin des simulations (voire même dès le début)
- ▶ Donc on peut exploiter ces probabilités pour améliorer les simulations suivantes (les rendre plus réalistes) : on favorise les simulations qui passent par les noeuds avec une forte probabilité de gagner au lieu de choisir aléatoirement
- ▶ On introduit un critère permettant d'explorer les branches qui ne l'ont pas été suffisamment

Recherche d'arbre Monte Carlo vs Algorithme aléatoire uniforme



- ▶ L'algorithme aléatoire uniforme génère N simulations indépendantes par coup possible à partir de l'état initial (état courant de la partie)
- ▶ MCTS construit un arbre de plus en plus grand où chaque nœud est associé à une estimation de la probabilité de gagner et en se concentrant sur les coups qui paraissent meilleurs

Recherche d'arbre Monte Carlo

Algorithme UCT (*Upper Confidence bounds in Trees*)

Tant que (il reste du temps)

1. **Sélectionner** récursivement à partir de la racine le nœud avec la plus grande B -valeur jusqu'à un nœud terminal ou un avec un fils non développé
2. **Développer** le nœud sélectionné (si nécessaire) et sélectionner aléatoirement un fils parmi les fils non explorés (ou le nœud lui-même s'il est terminal)
3. **Simuler** la fin de la partie avec une marche aléatoire
4. **Mettre à jour** les B -valeurs de tous les nœuds sur le chemin de la racine au nœud sélectionné en remontant la récompense r de la position finale

B -valeur du nœud i

$$B(i) = \pm \mu(i) + c \sqrt{\frac{\ln N(\text{parent}(i))}{N(i)}}$$

- ▶ \pm : alternance des joueurs, $+$ si $\text{parent}(i)$ est un nœud Max, $-$ si Min
- ▶ $\mu(i) = \frac{1}{N(i)} \sum_{\mathbf{s}: i \in \mathbf{s}} r(\mathbf{s})$: moyenne de la récompense $r(\mathbf{s})$ sur l'ensemble des simulations \mathbf{s} passant par i
si $r() \in \{0, 1\}$ (perdu ou gagné), alors $\mu(i)$ estime la probabilité de gagner
- ▶ $N(i)$: nombres de simulations passant par i
- ▶ c : constante permettant de régler le compromis exploitation / exploration

Recherche d'arbre Monte Carlo

B-valeur du nœud i

$$B(i) = \pm\mu(i) + c\sqrt{\frac{\ln N(\text{parent}(i))}{N(i)}}$$

Compromis Exploitation/Exploration

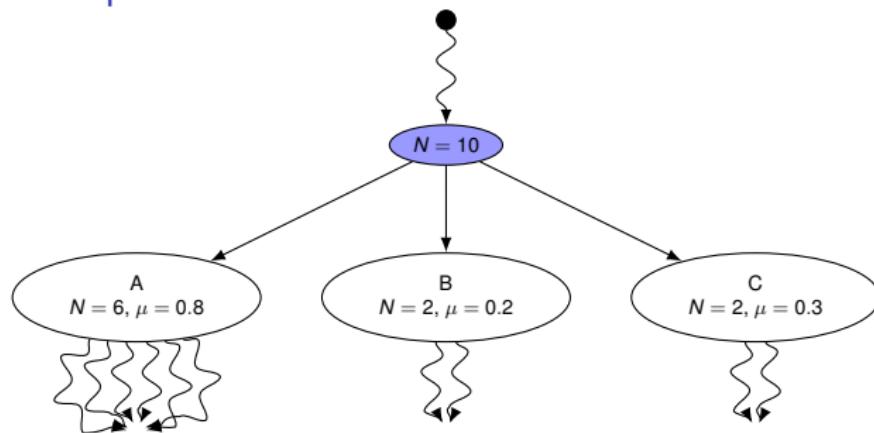
- ▶ Prendre $\mu(i)$ en compte : **exploitation** de l'information gagnée des précédentes simulations
- ▶ Le terme $\sqrt{\frac{\ln N(\text{parent}(i))}{N(i)}}$ favorise l'**exploration** :
 - ▶ il augmente pour les nœuds qui ne sont pas assez souvent sélectionnés : $N(i)$ est constant et $N(\text{parent}(i))$ augmente à chaque simulation dans laquelle on aurait pu choisir i
 - ▶ il décroît pour le nœud sélectionné : $N(i)$ augmente plus rapidement que $\ln N(\text{parent}(i))$
- ▶ Si un nœud paraît bon sur la base de $\mu(i)$, l'algorithme le sélectionnera mais il finira aussi par sélectionner un autre nœud après de nombreuses sélections concentrées sur le même nœud

Recherche d'arbre Monte Carlo – UCT

Exemple d'itération

Max

Min



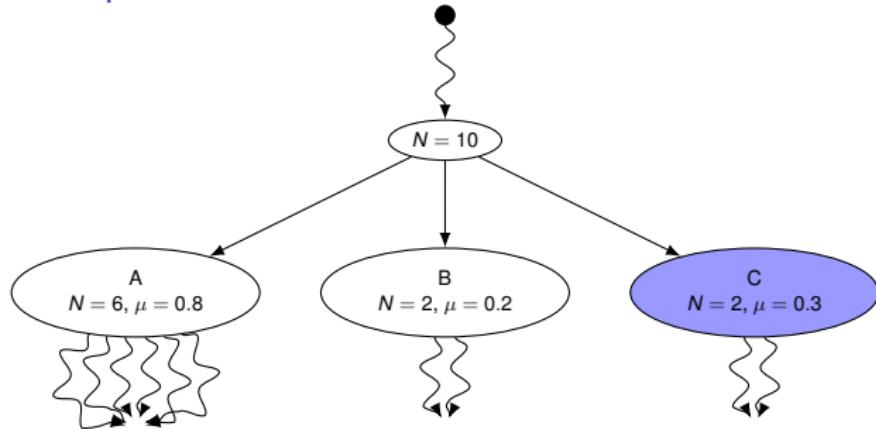
1. Sélection de A, B ou C

Recherche d'arbre Monte Carlo – UCT

Exemple d'itération

Max

Min



1. Sélection du nœud C (avec $c = 2$) :

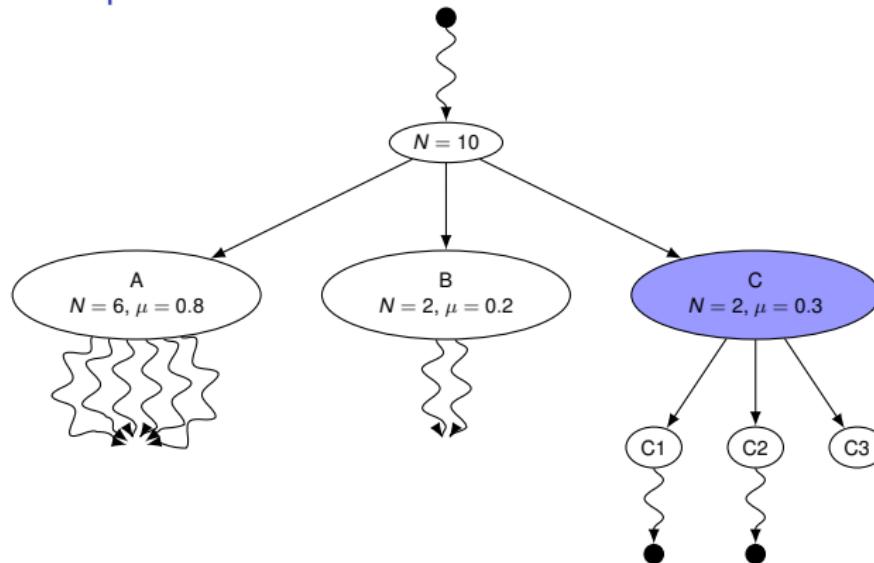
$$B(A) = 0.8 + 2\sqrt{\frac{\ln 10}{6}} \approx 2.04, \quad B(B) \approx 2.35, \quad B(C) \approx 2.45$$

Recherche d'arbre Monte Carlo – UCT

Exemple d'itération

Max

Min



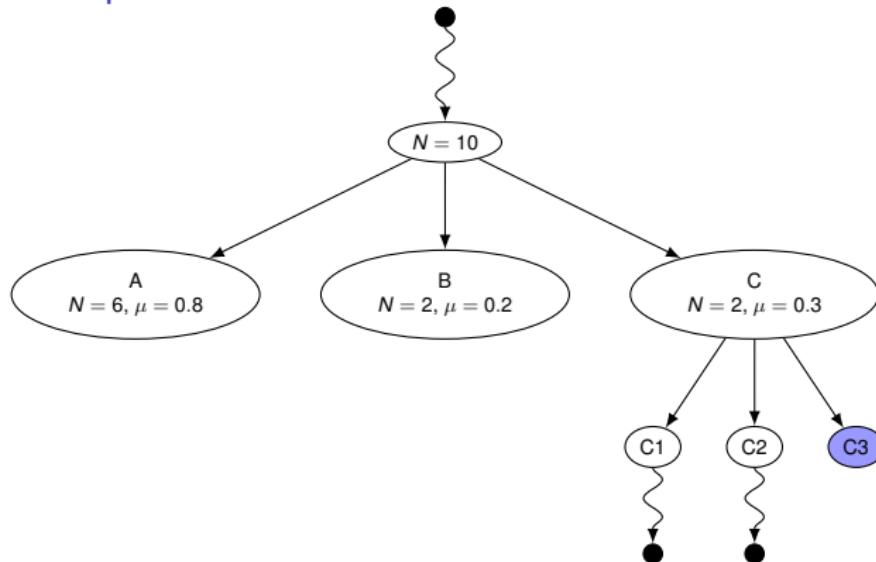
- Sélection** du nœud C
- Développement** du nœud sélectionné (déjà fait)

Recherche d'arbre Monte Carlo – UCT

Exemple d'itération

Max

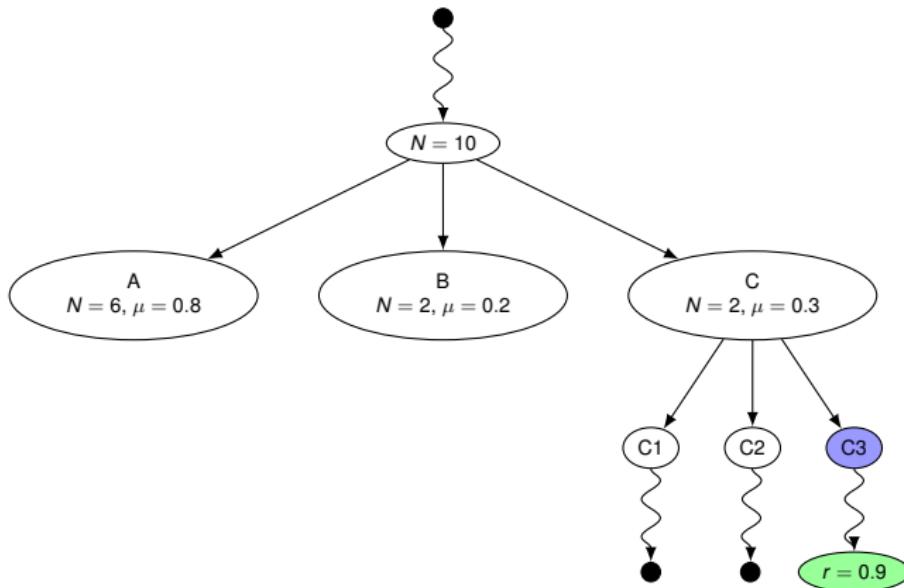
Min



1. **Sélection** du nœud C
2. **Développement** du nœud sélectionné (déjà fait) et sélection du fils non exploré C3

Recherche d'arbre Monte Carlo – UCT

Max

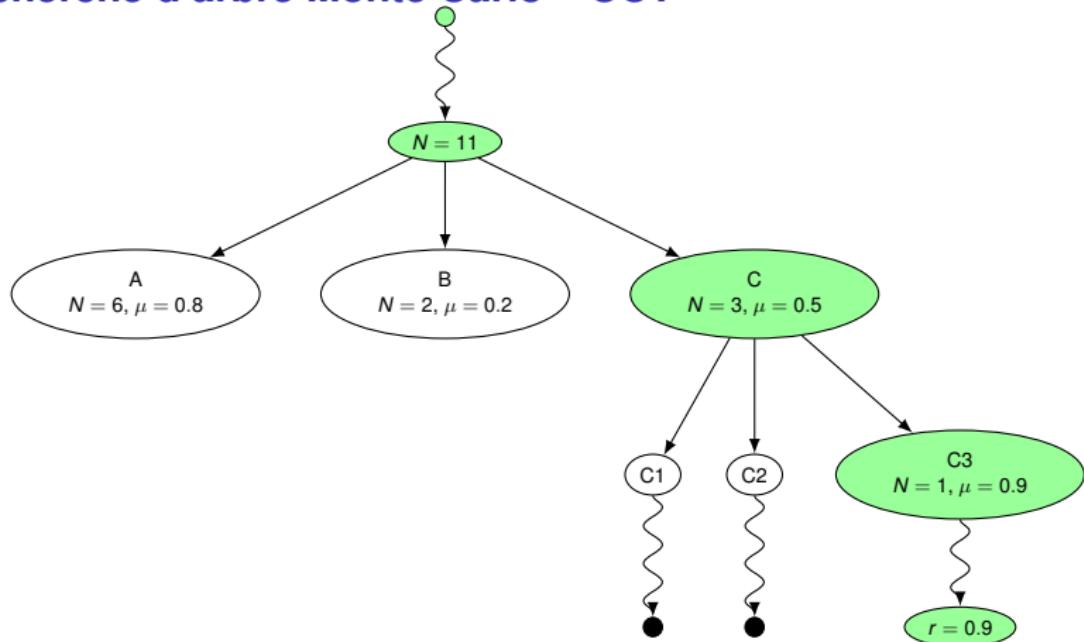


1. **Sélection** du nœud C
2. **Développement** du nœud sélectionné (déjà fait) et sélection du fils non exploré C3
3. **Simulation** aléatoire de la fin de la partie pour obtenir la récompense r

Recherche d'arbre Monte Carlo – UCT

Max

Min

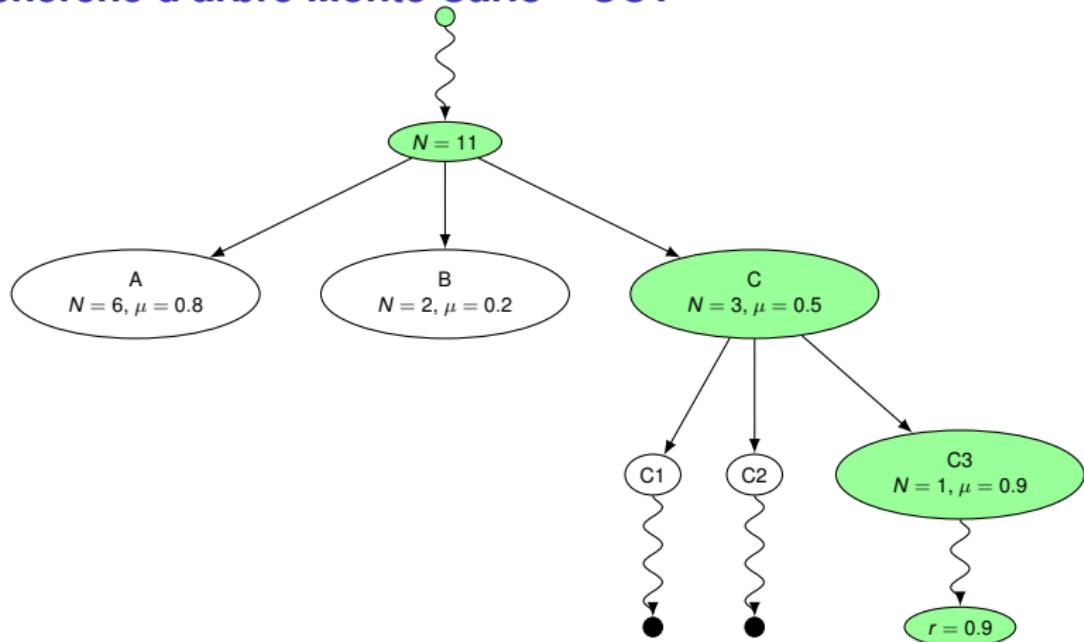


1. Sélection du nœud C
2. Développement du nœud sélectionné (déjà fait) et sélection du fils non exploré C3
3. Simulation aléatoire de la fin de la partie pour obtenir la récompense r
4. Mise à jour des N et μ sur tout le chemin de C3 à la racine

Recherche d'arbre Monte Carlo – UCT

Max

Min



1. Sélection du nœud C
2. Développement du nœud sélectionné (déjà fait) et sélection du fils non exploré C3
3. Simulation aléatoire de la fin de la partie pour obtenir la récompense r
4. Mise à jour des N et μ sur tout le chemin de C3 à la racine

Remarque : si $r \in \{0, 1\}$, $\mu(i) = \frac{\text{nb de victoires passant par } i}{N}$

Recherche d'arbre Monte Carlo

Mise à jour des B -valeurs en pratique

$$B(i) = \pm \mu(i) + c \sqrt{\frac{\ln N(\text{parent}(i))}{N(i)}} \quad \text{avec} \quad \mu(i) = \frac{1}{N(i)} \sum_{\mathbf{s}: i \in \mathbf{s}} r(\mathbf{s})$$

- ▶ Mémoriser $\mu(i)$ et $N(i)$ au lieu de $B(i)$
- ▶ Mise à jour :

$$\mu(i) \leftarrow \frac{N(i)\mu(i) + r}{N(i) + 1} \quad ; \quad N(i) \leftarrow N(i) + 1$$

ou alors :

- ▶ Mémoriser $S(i) = \sum_{\mathbf{s}: i \in \mathbf{s}} r(\mathbf{s})$ et calculer $\mu(i) = \frac{S(i)}{N(i)}$
- ▶ Mise à jour plus simple : $S(i) \leftarrow S(i) + r$
- ▶ Si la récompense est binaire (0 : perdu ou 1 : gagné), alors $S(i) = \text{nb de victoires passant par } i$

Recherche d'arbre Monte Carlo

Mise à jour des B -valeurs en pratique

$$B(i) = \pm\mu(i) + c\sqrt{\frac{\ln N(\text{parent}(i))}{N(i)}} \quad \text{avec} \quad \mu(i) = \frac{1}{N(i)} \sum_{\mathbf{s}: i \in \mathbf{s}} r(\mathbf{s})$$

- ▶ Mémoriser $\mu(i)$ et $N(i)$ au lieu de $B(i)$
- ▶ Mise à jour :

$$\mu(i) \leftarrow \frac{N(i)\mu(i) + r}{N(i) + 1} \quad ; \quad N(i) \leftarrow N(i) + 1$$

ou alors :

- ▶ Mémoriser $S(i) = \sum_{\mathbf{s}: i \in \mathbf{s}} r(\mathbf{s})$ et calculer $\mu(i) = \frac{S(i)}{N(i)}$
- ▶ Mise à jour plus simple : $S(i) \leftarrow S(i) + r$
- ▶ Si la récompense est binaire (0 : perdu ou 1 : gagné), alors $S(i) = \text{nb de victoires passant par } i$

Remarque : le signe codant l'alternance entre Max et Mini peut être pris en compte soit dans $B(i) = \pm\mu(i) + \dots$, soit dans $S(i) = \sum_{\mathbf{s}: i \in \mathbf{s}} \pm r(\mathbf{s})$

Recherche d'arbre Monte Carlo

Choix du coup à jouer

Tant que (il reste du temps)

...

Retourner le coup qui maximise

- ▶ soit la récompense attendue $\mu(i)$ (règle "max")
- ▶ soit le nombre de simulations $N(i)$ (règle "robuste")
- ▶ soit la récompense attendue $\mu(i)$ et le nombre de simulations (règle "max-robuste") (si aucun coup ne maximise les deux, on continue à chercher)

Améliorations

- ▶ Biasser la sélection et/ou les simulations
 - par ex. en choisissant toujours les coups gagnants ou qui empêchent l'autre joueur de gagner en un coup ; ou en repérant des motifs particuliers
- ▶ Utiliser une fonction d'évaluation pour limiter la profondeur des simulations
- ▶ Si l'ordre des coups n'a pas d'importance : mettre à jour les valeurs de tous les nœuds correspondant aux actions choisies pendant la simulation
- ▶ Beaucoup d'autres techniques plus ou moins spécifiques à chaque jeu

Succès de l'IA dans les jeux

IBM DeepBlue

- ▶ Aux échecs, $b \approx 30$, $n \approx 40$
- ▶ Au début et à la fin : répertoires d'ouvertures et de fin de partie
- ▶ Au milieu : alpha-beta
- ▶ Fonction d'évaluation avec +8000 paramètres + machine parallèle dédiée

AlphaGO

- ▶ Au GO, $b \approx 360$, $n \approx 400$, pas de fonction d'évaluation connue
- ▶ La machine doit apprendre la fonction d'évaluation
 1. Apprentissage d'un modèle imitateur (réseau de neurones profond) à partir d'une collection de parties jouées par des humains
 2. Le modèle joue contre lui-même pour s'améliorer et générer des millions de parties
 3. Apprentissage de la fonction d'évaluation à partir de ces parties
 4. MCTS biaisé par la fonction d'évaluation + simulations orientées par le modèle imitateur
+ beaucoup de CPU et GPU en parallèle

Recherche d'arbre Monte Carlo UCT et bandits manchots

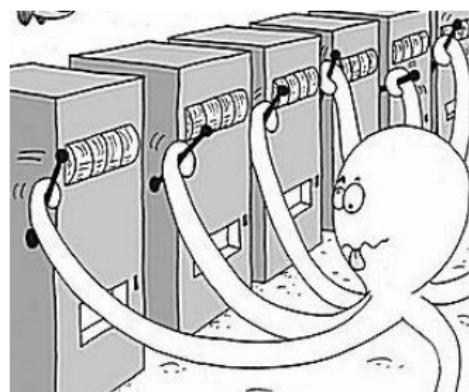
Recherche d'arbre Monte Carlo UCT et bandits manchots



Recherche d'arbre Monte Carlo UCT et bandits manchots

Bandits manchots à K bras (machines à sous)

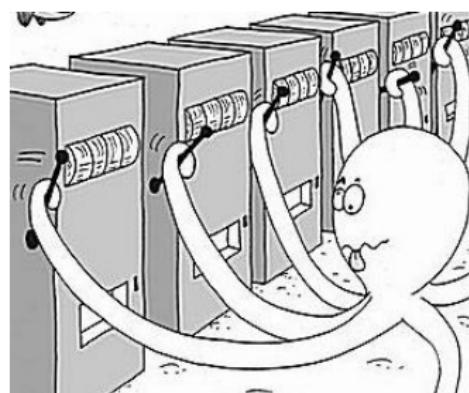
- ▶ Chaque bras a une probabilité **inconnue** de gagner
- ▶ A chaque étape, l'algorithme choisit un bras à tirer
- ▶ BUT : maximiser les gains en essayant de tirer le plus souvent possible le meilleur bras
- ▶ **Compromis Exploration / Exploitation** : tirer le bras que l'on pense être le meilleur (exploitation) tout en tirant suffisamment les autres pour améliorer notre connaissance des probabilités de gagner (exploration)



Recherche d'arbre Monte Carlo UCT et bandits manchots

Bandits manchots à K bras (machines à sous)

- ▶ Chaque bras a une probabilité **inconnue** de gagner
- ▶ A chaque étape, l'algorithme choisit un bras à tirer
- ▶ BUT : maximiser les gains en essayant de tirer le plus souvent possible le meilleur bras
- ▶ **Compromis Exploration / Exploitation** : tirer le bras que l'on pense être le meilleur (exploitation) tout en tirant suffisamment les autres pour améliorer notre connaissance des probabilités de gagner (exploration)



MCTS-UCT considère chaque décision à prendre à chaque nœud comme un problème de bandit manchot

Rappels de probabilités

Espace probabilisé (Ω, \mathcal{A}, P)

- ▶ **Univers** Ω : ensemble des résultats d'expérience possibles
- ▶ **Tribu (σ -algèbre) \mathcal{A} de Ω** : ensemble de tous les événements
 - tribu de $\Omega = ensemble de sous-ensembles de $\Omega tel que i) $\Omega \in \mathcal{A}, ii) $A \in \mathcal{A} \Rightarrow \bar{A} = (\Omega \setminus A) \in \mathcal{A}, iii) $A, B \in \mathcal{A} \Rightarrow A \cup B \in \mathcal{A}, iv) si $A_i \in \mathcal{A} pour tout $i \in \mathbb{N} alors $\cup_{i \in \mathbb{N}} A_i \in \mathcal{A}$$$$$$$$*
- ▶ **Événement A** : $A \subseteq \Omega$ et $A \in \mathcal{A}$
- ▶ **Mesure de probabilité P** : fonction associant une probabilité entre 0 et 1 à un événement

$$P : \mathcal{A} \rightarrow [0, 1]$$

telle que

- ▶ $P(\Omega) = 1$ et $P(\emptyset) = 0$
- ▶ Pour un nombre fini d'événements A_i disjoints : $P(\bigcup_i A_i) = \sum_i P(A_i)$

Rappels de probabilités

Espace probabilisé (Ω, \mathcal{A}, P)

- ▶ **Univers** Ω : ensemble des résultats d'expérience possibles
- ▶ **Tribu (σ -algèbre) \mathcal{A} de Ω** : ensemble de tous les événements
 - tribu de $\Omega = ensemble de sous-ensembles de $\Omega tel que i) $\Omega \in \mathcal{A}, ii) $A \in \mathcal{A} \Rightarrow \bar{A} = (\Omega \setminus A) \in \mathcal{A}, iii) $A, B \in \mathcal{A} \Rightarrow A \cup B \in \mathcal{A}, iv) si $A_i \in \mathcal{A} pour tout $i \in \mathbb{N} alors $\cup_{i \in \mathbb{N}} A_i \in \mathcal{A}$$$$$$$$*
- ▶ **Événement A** : $A \subseteq \Omega$ et $A \in \mathcal{A}$
- ▶ **Mesure de probabilité P** : fonction associant une probabilité entre 0 et 1 à un événement

$$P : \mathcal{A} \rightarrow [0, 1]$$

telle que

- ▶ $P(\Omega) = 1$ et $P(\emptyset) = 0$
- ▶ Pour un nombre fini d'événements A_i disjoints : $P(\bigcup_i A_i) = \sum_i P(A_i)$

Propriétés d'une probabilité

- ▶ $P(\bar{A}) = 1 - P(A)$ $\bar{A} = \Omega \setminus A$
- ▶ $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- ▶ Si A et B disjoints : $P(A \cup B) = P(A) + P(B)$

Rappels de probabilités

Exemple : le lancé d'un dé

- ▶ $\Omega = \{\square, \square\cdot, \square\cdot\cdot, \square\cdot\cdot\cdot, \square\cdot\cdot\cdot\cdot, \square\cdot\cdot\cdot\cdot\cdot\}$
- ▶ $\mathcal{A} = \mathcal{P}(\Omega) = \{\{\square\}, \{\square\cdot\}, \dots, \{\square, \square\cdot\}, \dots, \{\square, \square\cdot, \square\cdot\cdot\}, \dots, \dots\}$
 \mathcal{A} est l'ensemble de toutes les sous-parties de Ω
- ▶ 1 événement : $\{\square, \square\cdot\}$, 1 autre événement : $\{\square\cdot\}$
- ▶ Probabilité : $P(\{\omega\}) = 1/6$ pour tout $\omega \in \Omega$
- ▶ Union : $P(\{\square, \square\cdot\}) = P(\{\square\} \cup \{\square\cdot\}) = P(\{\square\}) + P(\{\square\cdot\}) = 2/6$
- ▶ Complément : $P(\{\square, \square\cdot, \square\cdot\cdot, \square\cdot\cdot\cdot, \square\cdot\cdot\cdot\cdot, \square\cdot\cdot\cdot\cdot\cdot\}) = P(\Omega) - P(\{\square\cdot\cdot\cdot\cdot\cdot\}) = 1 - 1/6 = 5/6$

Rappels de probabilités

Variables aléatoires (v.a.)

- ▶ Variable aléatoire réelle X : **fonction** associant une valeur numérique au résultat d'une expérience

$$X : \Omega \rightarrow \mathbb{R}$$

- ▶ Pour $A \subseteq \mathbb{R}$, l'ensemble des résultats d'expérience conduisant à une valeur de la v.a. dans A est un événement :

$$\{\omega \in \Omega \mid X(\omega) \in A\} \in \mathcal{A}$$

- ▶ On notera simplement X pour $X(\omega)$ et $\{X \in A\}$ pour $\{\omega \in \Omega \mid X(\omega) \in A\}$

Rappels de probabilités

Variables aléatoires (v.a.)

- ▶ Variable aléatoire réelle X : **fonction** associant une valeur numérique au résultat d'une expérience

$$X : \Omega \rightarrow \mathbb{R}$$

- ▶ Pour $A \subseteq \mathbb{R}$, l'ensemble des résultats d'expérience conduisant à une valeur de la v.a. dans A est un événement :

$$\{\omega \in \Omega \mid X(\omega) \in A\} \in \mathcal{A}$$

- ▶ On notera simplement X pour $X(\omega)$ et $\{X \in A\}$ pour $\{\omega \in \Omega \mid X(\omega) \in A\}$

Loi d'une variable aléatoire

- ▶ On définit une mesure de probabilité P_X sur \mathbb{R} muni de la tribu \mathcal{B} par

$$P_X : \mathcal{B} \rightarrow [0, 1]$$

$$P_X(A) = P(X \in A)$$

- ▶ P_X est la loi de X
- ▶ P_X est la probabilité image de P par la fonction X
- ▶ En pratique : utilisation de l'espace probabilisé $(\mathbb{R}, \mathcal{B}, P_X)$

Rappels de probabilités

Exemple : tirage aléatoire d'un individu dans une population

- ▶ Ω : ensemble des individus de la population
(chaque individu est un résultat possible)
- ▶ X : taille de l'individu en cm
- ▶ \mathcal{B} : tribu borélienne de \mathbb{R}
ensemble de toutes les parties de \mathbb{R} engendrées par union, intersections et compléments d'événements du type ($X < x$)
- ▶ Un exemple d'événement pour l'espace probabilisé $(\mathbb{R}, \mathcal{B}, P_X)$:
[170, 180]
- ▶ Correspond à l'événement $\{\omega \in \Omega | X(\omega) \in [170, 180]\}$, noté de manière courte $\{X \in [170, 180]\}$, pour l'espace probabilisé (Ω, \mathcal{A}, P)
- ▶ Probabilité : $P_X([170, 180]) = P(X \in [170, 180]) = 0.3$
(probabilité de tiré aléatoirement un individu dont la mesure de la taille soit entre 170 et 180 cm)

Rappels de probabilités

Variables aléatoires discrètes

- Une v.a. discrète Y ne peut prendre qu'un nombre fini de valeurs distinctes

$$Y \in \mathcal{Y} = \{y_k\}_{1 \leq k \leq n}$$

- La loi P_Y d'une v.a. discrète est donnée par la somme des probabilités de chacune des valeurs possibles :

$$P_Y(A) = \sum_{y \in \mathcal{Y} \cap A} P(Y = y) = \sum_{y_k \in A} P(Y = y_k)$$

Variables aléatoires continues

- Une v.a. continue X prend ses valeurs dans \mathbb{R}
- Densité de probabilité $p : \mathbb{R} \rightarrow \mathbb{R}^+$
- X suit la loi de **densité** p si :

$$P_X(A) = P(X \in A) = \int_A p(x) dx$$

- Remarque :

$$\text{pour } a \in \mathbb{R}, \quad P_X(a) = P(X = a) = \int_a^a p(x) dx = 0$$

Lois usuelles

Exemples de lois discrètes

- Loi uniforme sur $\{1, \dots, n\}$:

$$\forall y \in \{1, \dots, n\}, \quad P(Y = y) = \frac{1}{n}$$

- Loi de Bernoulli de paramètre $b \in [0, 1]$:

$$P(Y = 1) = b, \quad \text{et} \quad P(Y = 0) = 1 - b$$

Exemples de densités de probabilité

- Densité de la loi uniforme sur $[a, b]$:

$$p(x) = \begin{cases} \frac{1}{b-a}, & \text{si } x \in [a, b] \\ 0, & \text{sinon} \end{cases}$$

- Densité de la loi gaussienne (ou normale) $\mathcal{N}(\mu, \sigma^2)$:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(x-\mu)^2}{2\sigma^2}\right)$$

Rappels de probabilités

Espérance : valeur "en moyenne" d'une variable aléatoire

Espérance d'une v.a. discrète

$$\mathbb{E}[Y] = \sum_{y \in \mathcal{Y}} y P(Y = y)$$

Espérance d'une v.a. continue à densité

$$\mathbb{E}[X] = \int_{\mathbb{R}} x p(x) dx$$

Rappels de probabilités

Espérance : valeur "en moyenne" d'une variable aléatoire

Espérance d'une v.a. discrète

$$\mathbb{E}[Y] = \sum_{y \in \mathcal{Y}} y P(Y = y)$$

$$\mathbb{E}[f(Y)] = \sum_{y \in \mathcal{Y}} f(y) P(Y = y)$$

Espérance d'une v.a. continue à densité

$$\mathbb{E}[X] = \int_{\mathbb{R}} x p(x) dx$$

$$\mathbb{E}[f(X)] = \int_{\mathbb{R}} f(x) p(x) dx$$

Rappels de probabilités

Espérance : valeur "en moyenne" d'une variable aléatoire

Espérance d'une v.a. discrète

$$\mathbb{E}[Y] = \sum_{y \in \mathcal{Y}} y P(Y = y)$$

$$\mathbb{E}[f(Y)] = \sum_{y \in \mathcal{Y}} f(y) P(Y = y)$$

Espérance d'une v.a. continue à densité

$$\mathbb{E}[X] = \int_{\mathbb{R}} x p(x) dx$$

$$\mathbb{E}[f(X)] = \int_{\mathbb{R}} f(x) p(x) dx$$

Linéarité de l'espérance

$$\mathbb{E}[aX] = a\mathbb{E}[X] \quad et \quad \mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

Rappels de probabilités

Fonction indicatrice \mathbb{I}

- ▶ $\mathbb{I}_X(A) = \mathbb{I}(X \in A) = 1$ si l'événement A est observé, 0 sinon
- ▶ Exemple : $\mathbb{I}(X > 5)$

Espérance de l'indicatrice d'un événement = Probabilité de l'événement

En discret :

$$\mathbb{E}_Y[\mathbb{I}_Y(A)] = \sum_{y \in \mathcal{Y}} \mathbb{I}_Y(A) P(Y = y) = \sum_{y \in \mathcal{Y} \cap A} P(Y = y) = P_Y(A)$$

En continu :

$$\mathbb{E}_X[\mathbb{I}_X(A)] = \int_{\mathbb{R}} \mathbb{I}_X(A) p(x) dx = \int_A 1 p(x) dx + \int_{\bar{A}} 0 p(x) dx = P_X(A)$$

Bandits manchots à K bras (machines à sous)

- ▶ v.a. X_i : gain du bras i , de loi **inconnue** P_i de moyenne $\mathbb{E}[X_i] = \mu(i)$
- ▶ Le meilleur bras i^* a la plus grande moyenne : $i^* = \arg \max_{i=1,\dots,K} \mu(i)$
- ▶ Un algorithme/une politique A : $t \mapsto i$ choisi un bras i à chaque étape t

Bandits manchots à K bras (machines à sous)

- ▶ v.a. X_i : gain du bras i , de loi **inconnue** P_i de moyenne $\mathbb{E}[X_i] = \mu(i)$
- ▶ Le meilleur bras i^* a la plus grande moyenne : $i^* = \arg \max_{i=1,\dots,K} \mu(i)$
- ▶ Un algorithme/une politique A : $t \mapsto i$ choisi un bras i à chaque étape t
- ▶ Si on tire T fois le bras i , on gagne

$$\sum_{t=1}^T X_{i,t} \quad X_{i,t} : \text{copie indépendante de } X_i$$

Donc on peut espérer gagner au mieux

$$Gain^*(T) = \mathbb{E} \left[\sum_{t=1}^T X_{i^*,t} \right] = \sum_{t=1}^T \mathbb{E}[X_{i^*,t}] = \sum_{t=1}^T \mu(i^*) = T\mu(i^*)$$

Si A est déterministe et indépendant des $X_{i,t}$, on espère gagner

$$Gain(T) = \mathbb{E} \left[\sum_{t=1}^T X_{A(t),t} \right] = \sum_{t=1}^T \mu(A(t))$$

Bandits manchots à K bras (machines à sous)

- ▶ v.a. X_i : gain du bras i , de loi **inconnue** P_i de moyenne $\mathbb{E}[X_i] = \mu(i)$
- ▶ Le meilleur bras i^* a la plus grande moyenne : $i^* = \arg \max_{i=1,\dots,K} \mu(i)$
- ▶ Un algorithme/une politique A : $t \mapsto i$ choisi un bras i à chaque étape t
- ▶ Si on tire T fois le bras i , on gagne

$$\sum_{t=1}^T X_{i,t} \quad X_{i,t} : \text{copie indépendante de } X_i$$

Donc on peut espérer gagner au mieux

$$Gain^*(T) = \mathbb{E} \left[\sum_{t=1}^T X_{i^*,t} \right] = \sum_{t=1}^T \mathbb{E}[X_{i^*,t}] = \sum_{t=1}^T \mu(i^*) = T\mu(i^*)$$

Si A est déterministe et indépendant des $X_{i,t}$, on espère gagner

$$Gain(T) = \mathbb{E} \left[\sum_{t=1}^T X_{A(t),t} \right] = \sum_{t=1}^T \mu(A(t))$$

- ▶ Le **regret** est la somme de ce qu'on aurait pu espérer gagner en plus :

$$Regret(T) = \mathbb{E} [Gain^*(T) - Gain(T)] = \mathbb{E} \left[\sum_{t=1}^T (\mu(i^*) - \mu(A(t))) \right]$$

On prend l'espérance car A dépend des résultats aléatoires précédents

Bandits manchots à K bras (machines à sous)

- ▶ Comment choisir quel bras tirer pour **minimiser le regret** ?
- ▶ Au début : n'importe lequel (aucune information permettant de choisir)
- ▶ Mais ensuite ? et les 100 coups suivants ? Toujours tirer celui qui gagne le plus souvent ? Essayer un autre ?

Bandits manchots à K bras (machines à sous)

- ▶ Comment choisir quel bras tirer pour **minimiser le regret** ?
- ▶ Au début : n'importe lequel (aucune information permettant de choisir)
- ▶ Mais ensuite ? et les 100 coups suivants ? Toujours tirer celui qui gagne le plus souvent ? Essayer un autre ?
- ▶ Le regret se reformule en fonction du nombre $N_T(i)$ d'essais de chaque bras parmi T et des différences $\Delta_i = \mu(i^*) - \mu(i)$:

$$\begin{aligned} \text{Regret}(T) &= \mathbb{E} \left[\sum_{t=1}^T (\mu(i^*) - \mu(A(t))) \right] \\ &= \mathbb{E} \left[\sum_{i=1}^K N_T(i) (\mu(i^*) - \mu(i)) \right] \\ &= \sum_{i=1}^K \mathbb{E}[N_T(i)] \Delta_i \end{aligned}$$

A devrait conduire à des petits $N_T(i)$ pour les grands Δ_i

Bandits manchots à K bras – Algorithmes gloutons

Algorithme glouton (*exploiter au maximum*)

- ▶ Jouer au moins une fois chaque bras pour observer les gains $x_{i,j}$
- ▶ Estimer $\mu(i)$ par $\hat{\mu}(i) = \frac{1}{N_t(i)} \sum_{j=1}^t x_{i,j} \mathbb{I}(A(j) = i)$
- ▶ Choisir ensuite le bras qui maximise $\hat{\mu}(i)$

Bandits manchots à K bras – Algorithmes gloutons

Algorithme glouton (*exploiter au maximum*)

- ▶ Jouer au moins une fois chaque bras pour observer les gains $x_{i,j}$
- ▶ Estimer $\mu(i)$ par $\hat{\mu}(i) = \frac{1}{N_t(i)} \sum_{j=1}^t x_{i,j} \mathbb{I}(A(j) = i)$
- ▶ Choisir ensuite le bras qui maximise $\hat{\mu}(i)$
- ▶ Risque de toujours choisir le même bras sous-optimal
⇒ Regret linéaire en T

Bandits manchots à K bras – Algorithmes gloutons

Algorithme glouton (*exploiter au maximum*)

- ▶ Jouer au moins une fois chaque bras pour observer les gains $x_{i,j}$
- ▶ Estimer $\mu(i)$ par $\hat{\mu}(i) = \frac{1}{N_t(i)} \sum_{j=1}^t x_{i,j} \mathbb{I}(A(j) = i)$
- ▶ Choisir ensuite le bras qui maximise $\hat{\mu}(i)$
- ▶ Risque de toujours choisir le même bras sous-optimal
⇒ Regret linéaire en T

Algorithme ϵ -glouton (*explorer à l'infini*)

- ▶ A chaque étape,
 - ▶ avec une probabilité $1 - \epsilon$, appliquer l'algorithme glouton
 - ▶ avec une probabilité ϵ , tirer un bras aléatoirement

Bandits manchots à K bras – Algorithmes gloutons

Algorithme glouton (*exploiter au maximum*)

- ▶ Jouer au moins une fois chaque bras pour observer les gains $x_{i,j}$
- ▶ Estimer $\mu(i)$ par $\hat{\mu}(i) = \frac{1}{N_t(i)} \sum_{j=1}^t x_{i,j} \mathbb{I}(A(j) = i)$
- ▶ Choisir ensuite le bras qui maximise $\hat{\mu}(i)$
- ▶ Risque de toujours choisir le même bras sous-optimal
⇒ Regret linéaire en T

Algorithme ϵ -glouton (*explorer à l'infini*)

- ▶ A chaque étape,
 - ▶ avec une probabilité $1 - \epsilon$, appliquer l'algorithme glouton
 - ▶ avec une probabilité ϵ , tirer un bras aléatoirement
- ▶ Mais à chaque étape, le regret instantané est

$$\begin{aligned}\mathbb{E}[\mu(i^*) - \mu(A(t))] &= \sum_{i=1}^K \mathbb{E}[(\mu(i^*) - \mu(i)) \mathbb{I}(A(t) = i)] = \sum_{i=1}^K \Delta_i P(A(t) = i) \\ &\geq \frac{\epsilon}{K} \sum_{i=1}^K \Delta_i = O(1)\end{aligned}$$

Bandits manchots à K bras – Algorithmes gloutons

Algorithme glouton (*exploiter au maximum*)

- ▶ Jouer au moins une fois chaque bras pour observer les gains $x_{i,j}$
- ▶ Estimer $\mu(i)$ par $\hat{\mu}(i) = \frac{1}{N_t(i)} \sum_{j=1}^t x_{i,j} \mathbb{I}(A(j) = i)$
- ▶ Choisir ensuite le bras qui maximise $\hat{\mu}(i)$
- ▶ Risque de toujours choisir le même bras sous-optimal
⇒ Regret linéaire en T

Algorithme ϵ -glouton (*explorer à l'infini*)

- ▶ A chaque étape,
 - ▶ avec une probabilité $1 - \epsilon$, appliquer l'algorithme glouton
 - ▶ avec une probabilité ϵ , tirer un bras aléatoirement
- ▶ Mais à chaque étape, le regret instantané est

$$\begin{aligned}\mathbb{E}[\mu(i^*) - \mu(A(t))] &= \sum_{i=1}^K \mathbb{E}[(\mu(i^*) - \mu(i)) \mathbb{I}(A(t) = i)] = \sum_{i=1}^K \Delta_i P(A(t) = i) \\ &\geq \frac{\epsilon}{K} \sum_{i=1}^K \Delta_i = O(1)\end{aligned}$$

- ▶ ⇒ Regret total linéaire en T

Bandits manchots à K bras – Algorithmes gloutons

Algorithme glouton (*exploiter au maximum*)

- ▶ Jouer au moins une fois chaque bras pour observer les gains $x_{i,j}$
- ▶ Estimer $\mu(i)$ par $\hat{\mu}(i) = \frac{1}{N_t(i)} \sum_{j=1}^t x_{i,j} \mathbb{I}(A(j) = i)$
- ▶ Choisir ensuite le bras qui maximise $\hat{\mu}(i)$
- ▶ Risque de toujours choisir le même bras sous-optimal
⇒ Regret linéaire en T

Algorithme ϵ -glouton (*explorer à l'infini*)

- ▶ A chaque étape,
 - ▶ avec une probabilité $1 - \epsilon$, appliquer l'algorithme glouton
 - ▶ avec une probabilité ϵ , tirer un bras aléatoirement
- ▶ Mais à chaque étape, le regret instantané est

$$\begin{aligned}\mathbb{E}[\mu(i^*) - \mu(A(t))] &= \sum_{i=1}^K \mathbb{E}[(\mu(i^*) - \mu(i)) \mathbb{I}(A(t) = i)] = \sum_{i=1}^K \Delta_i P(A(t) = i) \\ &\geq \frac{\epsilon}{K} \sum_{i=1}^K \Delta_i = O(1)\end{aligned}$$

- ▶ ⇒ Regret total linéaire en T

En faisant décroître ϵ avec t , on peut obtenir un regret logarithmique *à partir de la connaissance des Δ_i*

Bandits manchots à K bras – Algorithme UCB

Idée

- ▶ Imaginons que l'on dispose d'une estimation $\hat{\mu}(i)$ de $\mu(i)$ et d'un semi-intervalle de confiance $C(i)$ tel que :

$$\mu(i) \leq \hat{\mu}(i) + C(i) = B(i)$$

- ▶ Alors on pourrait choisir le bras qui maximise $B(i)$ pour
 - ▶ favoriser les bras que l'on croit bons (grand $\hat{\mu}(i)$)
 - ▶ sans oublier les bras qui *pourraient* être meilleurs (grand $C(i)$)

Bandits manchots à K bras – Algorithme UCB

Idée

- ▶ Imaginons que l'on dispose d'une estimation $\hat{\mu}(i)$ de $\mu(i)$ et d'un semi-intervalle de confiance $C(i)$ tel que :

$$\mu(i) \leq \hat{\mu}(i) + C(i) = B(i)$$

- ▶ Alors on pourrait choisir le bras qui maximise $B(i)$ pour
 - ▶ favoriser les bras que l'on croit bons (grand $\hat{\mu}(i)$)
 - ▶ sans oublier les bras qui *pourraient* être meilleurs (grand $C(i)$)

Algorithme UCB (*Upper Confidence Bound*)

Tirer chaque bras une fois pour $t \leq K$, puis le bras i qui maximise

$$B_t(i) = \hat{\mu}_{N_t(i)}(i) + \sqrt{\frac{3 \ln t}{2N_t(i)}}$$

- ▶ $\hat{\mu}_n(i) = \frac{1}{n} \sum_{j=1}^n x_{i,j}$: estimation de $\mu(i)$ à partir de n observations $x_{i,j}$ des gains du bras i
- ▶ $N_t(i)$: Nombre de tirages du bras i sur les t premiers coups

Bandits manchots à K bras – Algorithme UCB

Idée

- ▶ Imaginons que l'on dispose d'une estimation $\hat{\mu}(i)$ de $\mu(i)$ et d'un semi-intervalle de confiance $C(i)$ tel que :

$$\mu(i) \leq \hat{\mu}(i) + C(i) = B(i)$$

- ▶ Alors on pourrait choisir le bras qui maximise $B(i)$ pour
 - ▶ favoriser les bras que l'on croit bons (grand $\hat{\mu}(i)$)
 - ▶ sans oublier les bras qui *pourraient* être meilleurs (grand $C(i)$)

Algorithme UCB (*Upper Confidence Bound*)

Tirer chaque bras une fois pour $t \leq K$, puis le bras i qui maximise

$$B_t(i) = \hat{\mu}_{N_t(i)}(i) + \sqrt{\frac{3 \ln t}{2N_t(i)}}$$

- ▶ $\hat{\mu}_n(i) = \frac{1}{n} \sum_{j=1}^n x_{i,j}$: estimation de $\mu(i)$ à partir de n observations $x_{i,j}$ des gains du bras i
- ▶ $N_t(i)$: Nombre de tirages du bras i sur les t premiers coups

MCTS-UCT : appliquer UCB à chaque nœud de l'arbre

Explication de la formule de $B(i)$ (UCB)

Théorème (inégalité de Hoeffding)

Soit une suite (X_j) de n v.a. indépendantes et identiquement distribuées à valeurs dans $[0, 1]$ et de moyenne μ , alors

$$P \left\{ \mu > \frac{1}{n} \sum_{j=1}^n X_j + \epsilon \right\} \leq e^{-2n\epsilon^2}$$

Explication de la formule de $B(i)$ (UCB)

Théorème (inégalité de Hoeffding)

Soit une suite (X_j) de n v.a. indépendantes et identiquement distribuées à valeurs dans $[0, 1]$ et de moyenne μ , alors

$$P \left\{ \mu > \frac{1}{n} \sum_{j=1}^n X_j + \epsilon \right\} \leq e^{-2n\epsilon^2}$$

Pour UCB

On souhaite être de plus en plus sûr de ce que l'on dit, et donc on veut une probabilité de plus en plus faible que μ dépasse l'intervalle de confiance. Soit δ_t une borne que l'on se fixe sur cette probabilité à l'instant t , on pose

$$e^{-2n\epsilon_t^2} = \delta_t \quad \text{et par exemple} \quad \delta_t = t^{-3}$$

On en déduit le (semi)-intervalle de confiance

$$\epsilon_t = \sqrt{\frac{3 \ln t}{2n}}$$

Plus t est grand, plus on est sûr que $\mu \leq \hat{\mu} + \epsilon_t$. Cependant, l'intervalle de confiance ϵ_t augmente, mais de façon raisonnable (logarithmique)

Borne sur le regret de UCB

Théorème

Le regret de la politique UCB sur T tirages est logarithmique en T

$$\text{Regret}(T) \leq \left[\sum_{\Delta_i \neq 0} \frac{6}{\Delta_i} \right] \ln T + K \left(1 + \frac{\pi^2}{3} \right) = O(\ln T)$$

Borne sur le regret de UCB

Théorème

Le regret de la politique UCB sur T tirages est logarithmique en T

$$\text{Regret}(T) \leq \left[\sum_{\Delta_i \neq 0} \frac{6}{\Delta_i} \right] \ln T + K \left(1 + \frac{\pi^2}{3} \right) = O(\ln T)$$

- ▶ $\text{Regret}(T) = \sum_{i=1}^K \Delta_i \mathbb{E}[N_T(i)]$, donc il suffit de montrer que
 $\mathbb{E}[N_T(i)] \leq \frac{6 \ln T}{\Delta_i^2} + 1 + \frac{\pi^2}{3}$ pour $\Delta_i \neq 0$

Borne sur le regret de UCB

Théorème

Le regret de la politique UCB sur T tirages est logarithmique en T

$$\text{Regret}(T) \leq \left[\sum_{\Delta_i \neq 0} \frac{6}{\Delta_i} \right] \ln T + K \left(1 + \frac{\pi^2}{3} \right) = O(\ln T)$$

- ▶ $\text{Regret}(T) = \sum_{i=1}^K \Delta_i \mathbb{E}[N_T(i)]$, donc il suffit de montrer que
 $\mathbb{E}[N_T(i)] \leq \frac{6 \ln T}{\Delta_i^2} + 1 + \frac{\pi^2}{3}$ pour $\Delta_i \neq 0$
- ▶ Calculons le nombre de tirages du bras $i \neq i^*$: $N_T(i) = \sum_{t=1}^T \mathbb{I}(A(t) = i)$

Borne sur le regret de UCB

Théorème

Le regret de la politique UCB sur T tirages est logarithmique en T

$$\text{Regret}(T) \leq \left[\sum_{\Delta_i \neq 0} \frac{6}{\Delta_i} \right] \ln T + K \left(1 + \frac{\pi^2}{3} \right) = O(\ln T)$$

- ▶ $\text{Regret}(T) = \sum_{i=1}^K \Delta_i \mathbb{E}[N_T(i)]$, donc il suffit de montrer que
 $\mathbb{E}[N_T(i)] \leq \frac{6 \ln T}{\Delta_i^2} + 1 + \frac{\pi^2}{3}$ pour $\Delta_i \neq 0$
- ▶ Calculons le nombre de tirages du bras $i \neq i^*$: $N_T(i) = \sum_{t=1}^T \mathbb{I}(A(t) = i)$
 $N_T(i) \leq n + \text{nb de tirages après } n \text{ tirages}$

$$\begin{aligned} &\leq n + \sum_{t=n+1}^T \mathbb{I}(A(t) = i, N_{t-1}(i) \geq n) \\ &\leq n + \sum_{t=n+1}^T \underbrace{\mathbb{I}(\hat{\mu}_{t-1}(i^*) + C_{t-1, N_{t-1}(i^*)} \leq \hat{\mu}_{t-1}(i) + C_{t-1, N_{t-1}(i)}, N_{t-1}(i) \geq n)}_{\text{necessaire pour avoir } A(t) = i} \\ &\leq n + \sum_{t=n}^{T-1} \mathbb{I}(\hat{\mu}_t(i^*) + C_{t, N_t(i^*)} \leq \hat{\mu}_t(i) + C_{t, N_t(i)}, N_t(i) \geq n) \end{aligned}$$

Avec $C_{t,N} = \sqrt{\frac{3 \ln t}{2N}}$

Borne sur le regret de UCB I

- $B_t(i^*) = \hat{\mu}_{N_t(i^*)}(i^*) + C_{t,N_t(i^*)} \leq \hat{\mu}_{N_t(i)}(i) + C_{t,N_t(i)} = B_t(i)$ nécessite soit $\mu(i^*)$ très sous-estimé, soit $\mu(i)$ très sur-estimé, soit une petite différence Δ_i :

$$\hat{\mu}_{N_t(i^*)}(i^*) \leq \mu(i^*) - C_{t,N_t(i^*)} \quad (1)$$

$$\hat{\mu}_{N_t(i)}(i) \geq \mu(i) + C_{t,N_t(i)} \quad (2)$$

$$\Delta_i = \mu(i^*) - \mu(i) < 2C_{t,N_t(i)} \quad (3)$$

Par ex., si pas (1) ni (3), $\hat{\mu}_{N_t(i^*)}(i^*) + C_{t,N_t(i^*)} > \mu(i^*) \geq \mu(i) + 2C_{t,N_t(i)}$ et il faut $\hat{\mu}_{N_t(i)}(i) + C_{t,N_t(i)} \geq \mu(i) + 2C_{t,N_t(i)} \Leftrightarrow (2)$; ...

- (3) $\Rightarrow N_t(i) < 6 \ln t / \Delta_i^2 \leq 6 \ln T / \Delta_i^2$, donc avec $N_t(i) \geq n = 6 \ln T / \Delta_i^2 + 1$, il faut (1) ou (2) pour avoir $A(t+1) = i$
- Avec la linéarité de l'espérance et la propriété $\mathbb{E}[\mathbb{I}(E)] = P(E)$:

$$\begin{aligned}\mathbb{E}[N_T(i)] &\leq \mathbb{E} \left[n + \sum_{t=n}^{T-1} \mathbb{I}((1) \text{ ou } (2)) \right] \\ &\leq n + \sum_{t=n}^{T-1} \mathbb{E} [\mathbb{I}((1) \text{ ou } (2))] \\ &\leq n + \sum_{t=n}^{T-1} P((1) \text{ ou } (2))\end{aligned}$$

Borne sur le regret de UCB II

- "Union bound" :

$$P((1) \text{ ou } (2)) \leq P(1) + P(2)$$

$$P(1) \leq P\left(\bigcup_{N=1}^t \{(1), N_t(i^*) = N\}\right) = \sum_{N=1}^t P((1), N_t(i^*) = N)$$

- Inégalité de Hoeffding : $P((1), N_t(i^*) = N) \leq e^{-2NC_{t,N}^2}$

$$\Rightarrow P(1) \leq \sum_{N=1}^t e^{-2NC_{t,N}^2} = \sum_{N=1}^t e^{-3 \ln t} = \sum_{N=1}^t t^{-3} = t^{-2}$$

- De même, $P(2) \leq t^{-2}$, donc

$$\mathbb{E}[N_T(i)] \leq n + \sum_{t=n}^{T-1} P((1) \text{ or } (2)) \leq \frac{6 \ln T}{\Delta_i^2} + 1 + \sum_{t=n}^{T-1} 2t^{-2}$$

- La série peut être bornée par une constante (grâce au problème de Bâle) :

$$\sum_{t=n}^{T-1} 2t^{-2} \leq 2 \sum_{t=1}^{\infty} \frac{1}{t^2} = \frac{\pi^2}{3}$$

- Donc

$$\mathbb{E}[N_T(i)] \leq \frac{6 \ln T}{\Delta_i^2} + 1 + \frac{\pi^2}{3} = O\left(\frac{6 \ln T}{\Delta_i^2}\right) = O(\ln T)$$

Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

L'intelligence artificielle face à l'incertain

Deux grands cas de figure

L'incertain dont on ne pourra jamais être sûr du résultat mais que l'on connaît bien

- ▶ Exemple : on "sait" que le résultat d'un lancer de dé est incertain et on "sait" (on arrive facilement à se convaincre que cela est assez proche de la réalité) que le ☐ va sortir avec une probabilité 1/6
- ▶ On connaît avec "certitude" un modèle de l'incertain (basé sur des probabilités)
- ▶ Mais on ne pourra jamais prédire le résultat avec certitude

L'intelligence artificielle face à l'incertain

Deux grands cas de figure

L'incertain dont on ne pourra jamais être sûr du résultat mais que l'on connaît bien

- ▶ Exemple : on "sait" que le résultat d'un lancer de dé est incertain et on "sait" (on arrive facilement à se convaincre que cela est assez proche de la réalité) que le ☐ va sortir avec une probabilité 1/6
- ▶ On connaît avec "certitude" un modèle de l'incertain (basé sur des probabilités)
- ▶ Mais on ne pourra jamais prédire le résultat avec certitude

L'incertain dont le résultat est certain mais sur lequel on ne sait rien



- ▶ Exemple : ces images représentent des chiffres de manière certaine, mais on ne sait pas comment, on ne connaît rien sur la fonction f dans $f(\text{image}) = \text{chiffre}$
- ▶ On ne sait pas comment modéliser le résultat d'une expérience
- ▶ Mais on sait que résultat n'est pas (ou peu) aléatoire

Deux approches pour gérer l'incertain

Modéliser directement l'incertitude à partir de nos connaissances

- ▶ Avec des probabilités, des variables aléatoires et des lois de probabilité
- ▶ Quand plusieurs variables sont "incertaines" : avec des lois jointes, des réseaux bayésiens...

Deux approches pour gérer l'incertain

Modéliser directement l'incertitude à partir de nos connaissances

- ▶ Avec des probabilités, des variables aléatoires et des lois de probabilité
- ▶ Quand plusieurs variables sont "incertaines" : avec des lois jointes, des réseaux bayésiens...

Apprendre pour combler nos lacunes



- ▶ En fait, faire en sorte que la machine apprenne à notre place
- ▶ A partir d'observations / d'exemples (des images et des chiffres)
- ▶ Avec des méthodes d'apprentissage automatique / statistique
- ▶ Pour obtenir une bonne approximation de f dans $f(\text{image}) = \text{chiffre}$

Deux approches pour gérer l'incertain

Modéliser directement l'incertitude à partir de nos connaissances

- ▶ Avec des probabilités, des variables aléatoires et des lois de probabilité
- ▶ Quand plusieurs variables sont "incertaines" : avec des lois jointes, des réseaux bayésiens...

Apprendre pour combler nos lacunes



- ▶ En fait, faire en sorte que la machine apprenne à notre place
- ▶ A partir d'observations / d'exemples (des images et des chiffres)
- ▶ Avec des méthodes d'apprentissage automatique / statistique
- ▶ Pour obtenir une bonne approximation de f dans $f(\text{image}) = \text{chiffre}$

2 parties du reste du cours

Réseaux bayésiens

- ▶ Permettent de modéliser des phénomènes aléatoires complexes impliquant de nombreuses variables aléatoires
Une loi jointe sur n variables à m valeurs possibles s'exprime a priori avec $O(m^n)$ nombres
- ▶ Par une approche graphique mettant en valeur les dépendances entre variables et limitant la taille de la représentation

Rappels de probabilités

Couple de v.a. discrètes $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, $|\mathcal{X}| < \infty$, $|\mathcal{Y}| < \infty$

- Loi jointe du couple : pour $A \subseteq \mathcal{X} \times \mathcal{Y}$

$$P_{X,Y}(A) = P((X, Y) \in A) = \sum_{(x,y) \in A} P(X = x, Y = y)$$

- Espérance :

$$\mathbb{E}_{(X,Y)}[f(X, Y)] = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} f(x, y) P(X = x, Y = y)$$

Couple de v.a. continues $(X, Y) \in \mathcal{X} \times \mathcal{Y} \subseteq \mathbb{R}^2$

- Densité de la loi jointe : $p_{X,Y}(x, y)$
- Loi jointe du couple : pour $A \subseteq \mathcal{X} \times \mathcal{Y}$

$$P_{X,Y}(A) = P((X, Y) \in A) = \int_A p_{X,Y}(x, y) dx dy$$

- Espérance : $\mathbb{E}_{(X,Y)}[f(X, Y)] = \iint_{\mathbb{R}^2} f(x, y) p_{X,Y}(x, y) dx dy$

Rappels de probabilités

Probabilités conditionnelles

- ▶ Probabilité de A sachant B : $P(A|B) = \frac{P(A,B)}{P(B)}$ (si $P(B) \neq 0$)
- ▶ Pour des v.a. discrètes : $P(X = x | Y = y) = \frac{P(X=x, Y=y)}{P(Y=y)}$
- ▶ Pour des v.a. continues : $p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$

A retenir : $P(X = x | Y = y)$ est une loi de probabilité fonction de y

Rappels de probabilités

Probabilités conditionnelles

- ▶ Probabilité de A sachant B : $P(A|B) = \frac{P(A,B)}{P(B)}$ (si $P(B) \neq 0$)
- ▶ Pour des v.a. discrètes : $P(X = x | Y = y) = \frac{P(X=x, Y=y)}{P(Y=y)}$
- ▶ Pour des v.a. continues : $p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$

A retenir : $P(X = x | Y = y)$ est une loi de probabilité fonction de y

Factorisation d'une loi jointe

- ▶ En discret :

$$P(X = x, Y = y) = P(X = x | Y = y)P(Y = y) = P(Y = y | X = x)P(X = x)$$

- ▶ En continu : $p_{X,Y}(x, y) = p_{X|Y}(x|y)p_Y(y) = p_{Y|X}(y|x)p_X(x)$

Rappels de probabilités

Probabilités conditionnelles

- ▶ Probabilité de A sachant B : $P(A|B) = \frac{P(A,B)}{P(B)}$ (si $P(B) \neq 0$)
- ▶ Pour des v.a. discrètes : $P(X = x|Y = y) = \frac{P(X=x, Y=y)}{P(Y=y)}$
- ▶ Pour des v.a. continues : $p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$

A retenir : $P(X = x|Y = y)$ est une loi de probabilité fonction de y

Factorisation d'une loi jointe

- ▶ En discret :

$$P(X = x, Y = y) = P(X = x|Y = y)P(Y = y) = P(Y = y|X = x)P(X = x)$$

- ▶ En continu : $p_{X,Y}(x,y) = p_{X|Y}(x|y)p_Y(y) = p_{Y|X}(y|x)p_X(x)$

Lois marginales et probabilités totales

- ▶ En discret :

$$P(X = x) = \sum_{y \in \mathcal{Y}} P(X = x, Y = y) = \sum_{y \in \mathcal{Y}} P(X = x|Y = y)P(Y = y)$$

Rappels de probabilités

Indépendance et indépendance conditionnelle

- X et Y sont indépendantes si et seulement si

$$P(X = x, Y = y) = P(X = x)P(Y = y)$$

- X et Y sont **conditionnellement indép.** sachant Z si et seulement si

$$P(X = x, Y = y|Z = z) = P(X = x|Z = z)P(Y = y|Z = z)$$

- Si X et Y sont conditionnellement indép. sachant Z , on a aussi

$$P(X = x|Y = y, Z = z) = P(X = x|Z = z)$$

Car

$$\begin{aligned} P(X = x|Y = y, Z = z) &= \frac{P(X = x, Y = y, Z = z)}{P(Y = y, Z = z)} = \frac{P(X = x, Y = y|Z = z)P(Z = z)}{P(Y = y, Z = z)} \\ &= \frac{P(X = x, Y = y|Z = z)P(Z = z)}{P(Y = y|Z = z)P(Z = z)} \\ &= \frac{P(X = x, Y = y|Z = z)}{P(Y = y|Z = z)} = \frac{P(X = x|Z = z)P(Y = y|Z = z)}{P(Y = y|Z = z)} \\ &= P(X = x|Z = z) \end{aligned}$$

Réseaux bayésiens

Vecteurs aléatoires

- ▶ $\boldsymbol{X} \in \mathbb{R}^n$ est un vecteur aléatoire contenant n v.a. continues
- ▶ $\boldsymbol{X} \in \{0, 1\}^n$ est un vecteur aléatoire binaire contenant n v.a. binaires
- ▶ La loi de \boldsymbol{X} est la loi de jointe de ses composantes :

$$P(\boldsymbol{X} = \mathbf{x}) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

Représenter des lois jointes en grande dimension

- ▶ Pour des vecteurs aléatoires binaires, $P(\boldsymbol{X} = \mathbf{x})$ peut s'écrire comme un tableau de $2^n - 1$ lignes contenant $2^n - 1$ paramètres

| x_1 | x_2 | ... | x_n | $P(\boldsymbol{X} = (x_1, x_2, \dots, x_n))$ |
|-------|-------|-----|-------|--|
| 0 | 0 | | 0 | 0.1 |
| 0 | 0 | | 1 | 0.8 |
| : | | | | : |

Réseaux bayésiens

Vecteurs aléatoires

- ▶ $\mathbf{X} \in \mathbb{R}^n$ est un vecteur aléatoire contenant n v.a. continues
- ▶ $\mathbf{X} \in \{0, 1\}^n$ est un vecteur aléatoire binaire contenant n v.a. binaires
- ▶ La loi de \mathbf{X} est la loi de jointe de ses composantes :

$$P(\mathbf{X} = \mathbf{x}) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

Représenter des lois jointes en grande dimension

- ▶ Pour des vecteurs aléatoires binaires, $P(\mathbf{X} = \mathbf{x})$ peut s'écrire comme un tableau de $2^n - 1$ lignes contenant $2^n - 1$ paramètres

| x_1 | x_2 | ... | x_n | $P(\mathbf{X} = (x_1, x_2, \dots, x_n))$ |
|-------|-------|-----|-------|--|
| 0 | 0 | | 0 | 0.1 |
| 0 | 0 | | 1 | 0.8 |
| : | | | | : |

- ▶ Difficile à implémenter pour n grand
- ▶ Besoin d'inclure des hypothèses/constraints supplémentaires

Réseaux bayésiens

Représentation graphique d'une loi jointe

- ▶ Chaque nœud correspond à une variable aléatoire
- ▶ Les arcs orientés déterminent les (in)dépendances conditionnelles chaque v.a. est conditionnellement indépendante de ses prédécesseurs sachant ses parents :

$$P(X_i = x_i | \mathbf{X}_{1:i-1} = \mathbf{x}_{1:i-1}) = P(X_i = x_i | \mathbf{X}_{parents(X_i)} = \mathbf{x}_{parents(X_i)})$$

Réseaux bayésiens

Représentation graphique d'une loi jointe

- ▶ Chaque nœud correspond à une variable aléatoire
- ▶ Les arcs orientés déterminent les (in)dépendances conditionnelles chaque v.a. est conditionnellement indépendante de ses prédécesseurs sachant ses parents :

$$P(X_i = x_i | \mathbf{X}_{1:i-1} = \mathbf{x}_{1:i-1}) = P(X_i = x_i | \mathbf{X}_{\text{parents}(X_i)} = \mathbf{x}_{\text{parents}(X_i)})$$

Factorisation d'une loi jointe par un réseau bayésien

$$P(\mathbf{X} = \mathbf{x}) = \prod_{i=1}^n P(X_i = x_i | \mathbf{X}_{\text{parents}(X_i)} = \mathbf{x}_{\text{parents}(X_i)})$$

Si les X_i sont ordonnées de telle sorte que $j \in \text{parents}(X_i) \Rightarrow j < i$,

$$\begin{aligned} P(\mathbf{X} = \mathbf{x}) &= P(X_n = x_n | \mathbf{X}_{1:n-1} = \mathbf{x}_{1:n-1}) P(\mathbf{X}_{1:n-1} = \mathbf{x}_{1:n-1}) \\ &= P(X_n = x_n | \mathbf{X}_{1:n-1} = \mathbf{x}_{1:n-1}) P(X_{n-1} = x_{n-1} | \mathbf{X}_{1:n-2} = \mathbf{x}_{1:n-2}) P(\mathbf{X}_{1:n-2} = \mathbf{x}_{1:n-2}) \\ &= \dots \\ &= \prod_{i=1}^n P(X_i = x_i | \mathbf{X}_{1:i-1} = \mathbf{x}_{1:i-1}) = \prod_{i=1}^n P(X_i = x_i | \mathbf{X}_{\text{parents}(X_i)} = \mathbf{x}_{\text{parents}(X_i)}) \end{aligned}$$

avec $1 : 0 = \emptyset$ et $\text{parents}(X_1) = \emptyset$

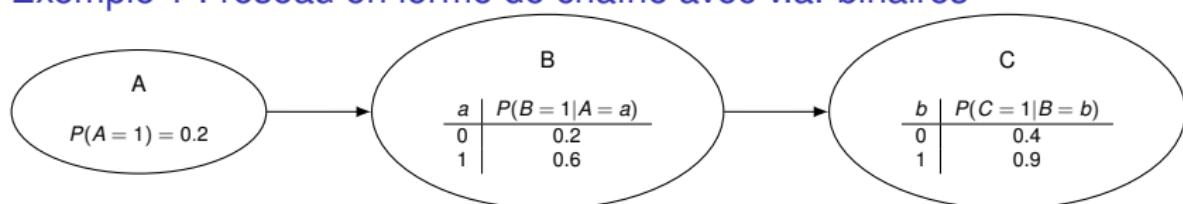
Réseaux bayésiens

Table de probabilités conditionnelles

- ▶ Chaque nœud correspond à une variable aléatoire X_i
- ▶ On associe une table de probabilités conditionnelles à chaque nœud pour définir

$$P(X_i = x_i \mid \mathbf{X}_{\text{parents}(X_i)} = \mathbf{x}_{\text{parents}(X_i)})$$

Exemple 1 : réseau en forme de chaîne avec v.a. binaires

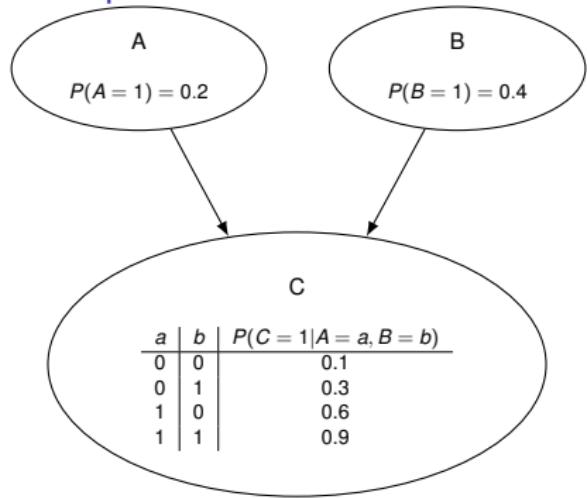


- ▶ 5 nombres suffisent ici pour coder la loi jointe de 3 variables binaires au lieu de $2^3 - 1 = 7$

$$\begin{aligned}P(A = 0, B = 1, C = 1) &= P(C = 1 | B = 1)P(B = 1 | A = 0)P(A = 0) \\&= 0.9 \times 0.2 \times 0.8 = 0.144\end{aligned}$$

Réseaux bayésiens – exemples avec des v.a. binaires

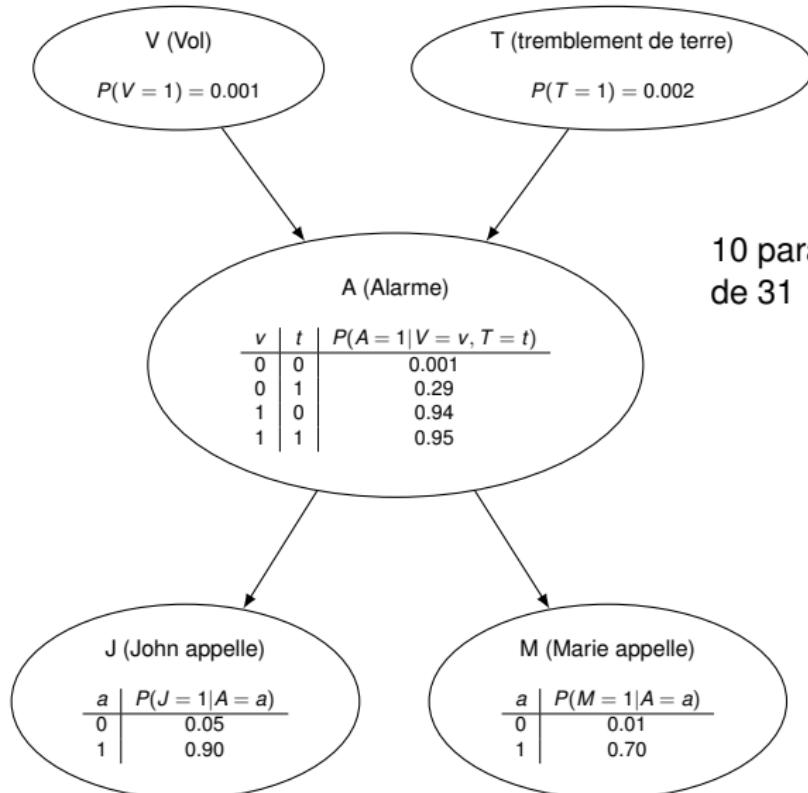
Exemple 2 : réseau en forme de V



- ▶ 6 nombres suffisent ici pour coder la loi jointe de 3 variables binaires au lieu de $2^3 - 1 = 7$

$$\begin{aligned}P(A = 0, B = 1, C = 1) &= P(C = 1 | A = 0, B = 1)P(A = 0)P(B = 1) \\&= 0.3 \times 0.8 \times 0.4 = 0.288\end{aligned}$$

Réseaux bayésiens



10 paramètres au lieu
de 31

$$\begin{aligned} & P(V = 1, T = 0, A = 1, J = 1, M = 0) \\ &= P(V = 1)P(T = 0)P(A = 1 | V = 1, T = 0)P(J = 1 | A = 1)P(M = 0 | A = 1) \end{aligned}$$

Réseaux bayésiens

Construction d'un réseaux

- ▶ Plusieurs topologies possibles pour la même loi jointe
- ▶ Certaines plus complexes et moins intuitives
- ▶ Dans un réseau simple, les arcs modélisent les liens de **cause à effet directs**
- ▶ Commencer le réseaux avec les causes premières (les événements qui semblent indépendants du reste)
- ▶ Ajouter les nœuds un par un dans l'ordre des causes à effets et en créant un jeu minimal d'arcs

Réseaux bayésiens

Construction d'un réseaux

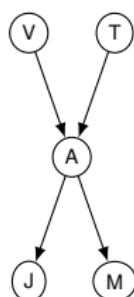
- ▶ Plusieurs topologies possibles pour la même loi jointe
- ▶ Certaines plus complexes et moins intuitives
- ▶ Dans un réseau simple, les arcs modélisent les liens de **cause à effet directs**
- ▶ Commencer le réseaux avec les causes premières (les événements qui semblent indépendants du reste)
- ▶ Ajouter les nœuds un par un dans l'ordre des causes à effets et en créant un jeu minimal d'arcs

Procédure itérative

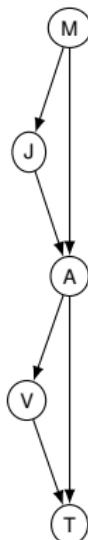
1. Choisir l'ensemble des variables et les ordonner : X_1, X_2, \dots, X_n
2. Pour $i = 1, \dots, n$
 - 2.1 ajouter un nœud pour la variable X_i
 - 2.2 ajouter des arcs en provenance des *parents*(X_i) déjà présents dans le **réseau** en respectant les hypothèses d'indépendances conditionnelles
 $P(X_i = x_i | \mathbf{X}_{\text{parents}(X_i)} = \mathbf{x}_{\text{parents}(x_i)}) = P(X_i = x_i | \mathbf{X}_{1:i-1} = \mathbf{x}_{1:i-1})$
 - 2.3 remplir la table des probabilités conditionnelles de X_i
valeurs de $P(X_i = x_i | \mathbf{X}_{\text{parents}(X_i)} = \mathbf{x}_{\text{parents}(x_i)})$ pour tout $\mathbf{x}_{\text{parents}(x_i)}$

Réseaux bayésiens

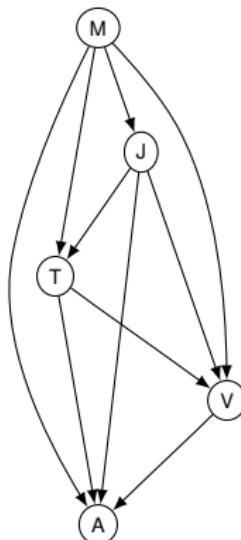
Exemple de topologies obtenues en ajoutant les nœuds dans un ordre différent



V, T, A, J, M
10 paramètres



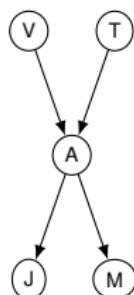
M, J, A, V, T
?? paramètres



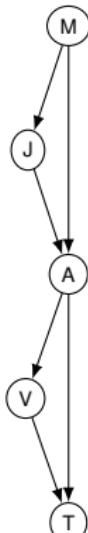
M, J, T, V, A
?? paramètres

Réseaux bayésiens

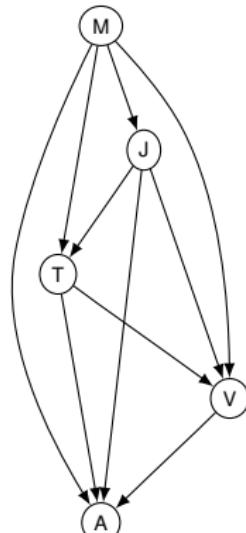
Exemple de topologies obtenues en ajoutant les nœuds dans un ordre différent



V, T, A, J, M
10 paramètres



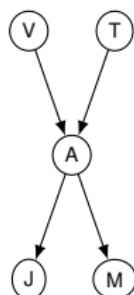
M, J, A, V, T
13 paramètres



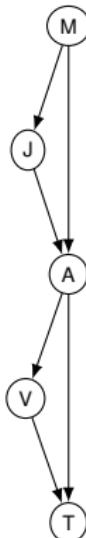
M, J, T, V, A
?? paramètres

Réseaux bayésiens

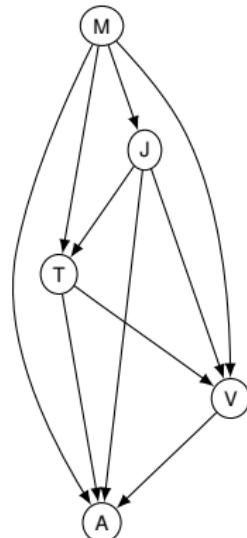
Exemple de topologies obtenues en ajoutant les nœuds dans un ordre différent



V, T, A, J, M
10 paramètres



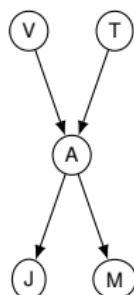
M, J, A, V, T
13 paramètres



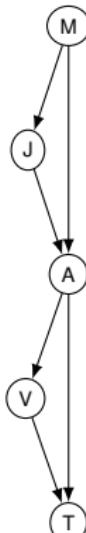
M, J, T, V, A
31 paramètres

Réseaux bayésiens

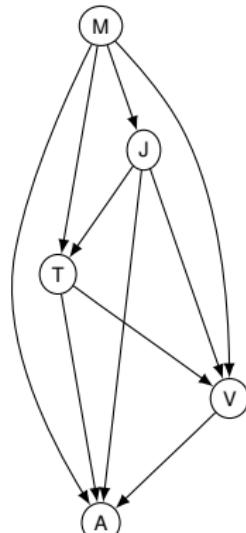
Exemple de topologies obtenues en ajoutant les nœuds dans un ordre différent



V, T, A, J, M
10 paramètres



M, J, A, V, T
13 paramètres



M, J, T, V, A
31 paramètres ($= 2^5 - 1$)

Réseaux bayésiens

D-séparation et indépendances conditionnelles

Notations

- ▶ $A \perp B$: A et B sont indépendantes
- ▶ $A \perp B | C$: A et B sont conditionnellement indépendantes sachant C

Réseaux bayésiens

D-séparation et indépendances conditionnelles

Notations

- ▶ $A \perp B$: A et B sont indépendantes
- ▶ $A \perp B | C$: A et B sont conditionnellement indépendantes sachant C

- ▶ Evidence $\mathbf{E} = \mathbf{e}$: ensemble de v.a. dont la valeur est connue / observée
- ▶ Comment savoir si $A \perp B | \mathbf{E}$?

D-séparation et indépendances conditionnelles

Notations

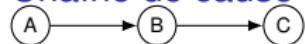
- ▶ $A \perp B$: A et B sont indépendantes
- ▶ $A \perp B | C$: A et B sont conditionnellement indépendantes sachant C

- ▶ Evidence $\mathbf{E} = \mathbf{e}$: ensemble de v.a. dont la valeur est connue / observée
- ▶ Comment savoir si $A \perp B | \mathbf{E}$?

- ▶ *d*-séparation : notion de séparation "directionnelle" ou "orientée"
- ▶ A et B sont *d*-séparées par $\mathbf{E} \Rightarrow A \perp B | \mathbf{E}$

Réseaux bayésiens – triplets élémentaires

Chaîne de cause à effet



- ▶ $A \perp B ?$
- ▶ $A \perp C ?$
- ▶ $A \perp C | B ?$

Réseaux bayésiens – triplets élémentaires

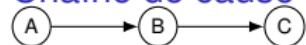
Chaîne de cause à effet



- ▶ $A \perp B$? Non
- ▶ $A \perp C$?
- ▶ $A \perp C | B$?

Réseaux bayésiens – triplets élémentaires

Chaîne de cause à effet



- ▶ $A \perp B$? Non
- ▶ $A \perp C$? Non
- ▶ $A \perp C | B$?

Réseaux bayésiens – triplets élémentaires

Chaîne de cause à effet



- ▶ $A \perp B$? Non
- ▶ $A \perp C$? Non
- ▶ $A \perp C | B$? Oui

$$\begin{aligned} P(A = a, C = c | B = b) &= \frac{P(A = a, B = b, C = c)}{P(B = b)} \\ &= \frac{P(A = a)P(B = b | A = a)P(C = c | B = b)}{P(B = b)} \\ &= P(A = a | B = b)P(C = c | B = b) \end{aligned}$$

Réseaux bayésiens – triplets élémentaires

Chaîne de cause à effet



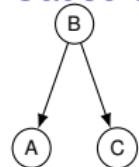
- ▶ $A \perp B$? Non
- ▶ $A \perp C$? Non
- ▶ $A \perp C | B$? Oui

$$\begin{aligned} P(A = a, C = c | B = b) &= \frac{P(A = a, B = b, C = c)}{P(B = b)} \\ &= \frac{P(A = a)P(B = b | A = a)P(C = c | B = b)}{P(B = b)} \\ &= P(A = a | B = b)P(C = c | B = b) \end{aligned}$$

- ▶ Après avoir observé B , A n'a plus d'influence sur ma croyance en C

Réseaux bayésiens – triplets élémentaires

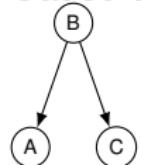
Cause commune



- ▶ $A \perp C$? Non
- ▶ $A \perp C|B$?

Réseaux bayésiens – triplets élémentaires

Cause commune

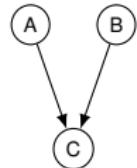


- ▶ $A \perp C$? Non
- ▶ $A \perp C|B$? Oui

$$\begin{aligned} P(A = a, C = c | B = b) &= \frac{P(A = a, B = b, C = c)}{P(B = b)} \\ &= \frac{P(A = a | B = b)P(B = b)P(C = c | B = b)}{P(B = b)} \\ &= P(A = a | B = b)P(C = c | B = b) \end{aligned}$$

Réseaux bayésiens – triplets élémentaires

Effet commun (structure en V)

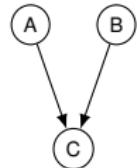


► $A \perp B ?$

► $A \perp B | C ?$

Réseaux bayésiens – triplets élémentaires

Effet commun (structure en V)



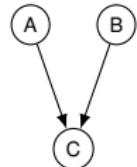
- ▶ $A \perp B$? Oui

$$\begin{aligned} P(A = a, B = b) &= \sum_c P(A = a, B = b, C = c) \\ &= P(A = a)P(B = b) \sum_c P(C = c | A = a, B = b) \\ &= P(A = a)P(B = b) \end{aligned}$$

- ▶ $A \perp B | C$?

Réseaux bayésiens – triplets élémentaires

Effet commun (structure en V)



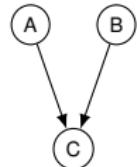
- ▶ $A \perp B$? Oui

$$\begin{aligned} P(A = a, B = b) &= \sum_c P(A = a, B = b, C = c) \\ &= P(A = a)P(B = b) \sum_c P(C = c | A = a, B = b) \\ &= P(A = a)P(B = b) \end{aligned}$$

- ▶ $A \perp B | C$? NON : après avoir observé C , A et B deviennent liées
Par exemple : $A =$ il fait beau ? $B =$ j'ai gagné au loto ? $C =$ je suis content ?

Réseaux bayésiens – triplets élémentaires

Effet commun (structure en V)



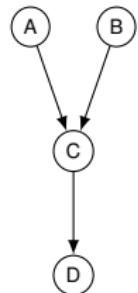
- ▶ $A \perp B$? Oui

$$\begin{aligned} P(A = a, B = b) &= \sum_c P(A = a, B = b, C = c) \\ &= P(A = a)P(B = b) \sum_c P(C = c | A = a, B = b) \\ &= P(A = a)P(B = b) \end{aligned}$$

- ▶ $A \perp B | C$? NON : après avoir observé C , A et B deviennent liées
Par exemple : A = il fait beau ? B = j'ai gagné au loto ? C = je suis content ?
- ▶ Utile pour le diagnostic
Par exemple : A = panne de batterie ? B = panne d'essence ? C = voiture ne démarre pas ?

Réseaux bayésiens – triplets élémentaires

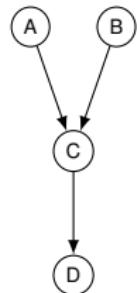
Effet commun (structure en V) observé "à distance"



- ▶ $A \perp B$? Oui
- ▶ $A \perp B|C$? Non
- ▶ $A \perp B|D$?

Réseaux bayésiens – triplets élémentaires

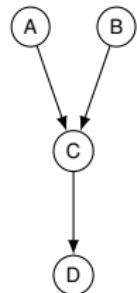
Effet commun (structure en V) observé "à distance"



- ▶ $A \perp B$? Oui
- ▶ $A \perp B|C$? Non
- ▶ $A \perp B|D$? Non : observer D me donne de l'information sur C qui me permet de lier A et B

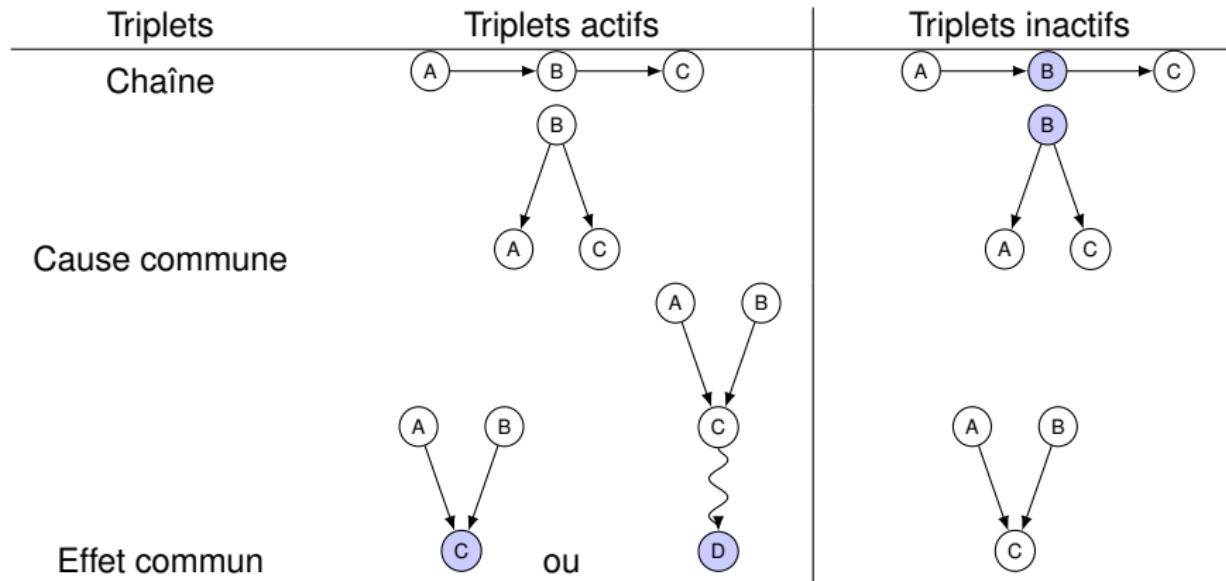
Réseaux bayésiens – triplets élémentaires

Effet commun (structure en V) observé "à distance"



- ▶ $A \perp B$? Oui
- ▶ $A \perp B|C$? Non
- ▶ $A \perp B|D$? Non : observer D me donne de l'information sur C qui me permet de lier A et B
- ▶ Exemples :
 - A = note < 6 en IA ? B = < 6 en LMC ? C = semestre non validé ? D = étudiant au rattrapage de Réseau ?
 - A = câble branché ? B = IP configurée ? C = connexion internet fonctionnelle ? D = page web s'affiche ?

Réseaux bayésiens – triplets élémentaires



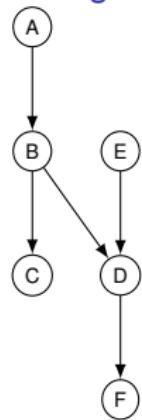
(B) : évidence (nœud observé)

Cas général : $X \perp Y | Z$?

- ▶ $X \perp Y | Z$ si et seulement si il n'existe pas de chemin (indépendamment de l'orientation des arcs) actif entre X et Y
- ▶ Chemin actif = chemin le long duquel tous les triplets sont actifs

Réseaux bayésiens

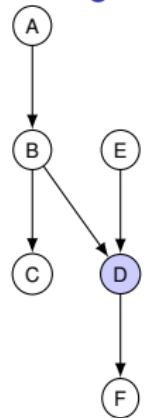
Cas général



► $A \perp F ?$

Réseaux bayésiens

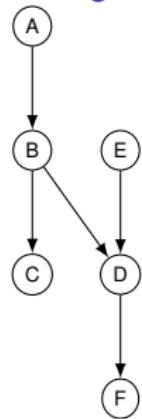
Cas général



- ▶ $A \perp\!\!\!\perp F$?
- ▶ $A \perp\!\!\!\perp F|D$?

Réseaux bayésiens

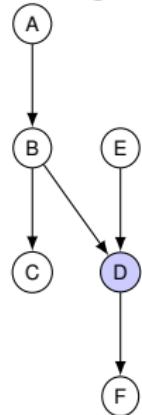
Cas général



- ▶ $A \perp F ?$
- ▶ $A \perp F|D ?$
- ▶ $A \perp E ?$

Réseaux bayésiens

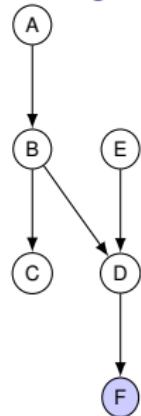
Cas général



- ▶ $A \perp F ?$
- ▶ $A \perp F|D ?$
- ▶ $A \perp E ?$
- ▶ $A \perp E|D ?$

Réseaux bayésiens

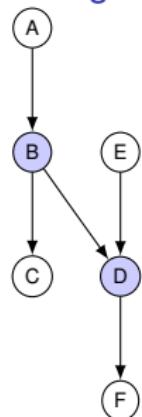
Cas général



- ▶ $A \perp\!\!\!\perp F ?$
- ▶ $A \perp\!\!\!\perp F|D ?$
- ▶ $A \perp\!\!\!\perp E ?$
- ▶ $A \perp\!\!\!\perp E|D ?$
- ▶ $A \perp\!\!\!\perp E|F ?$

Réseaux bayésiens

Cas général

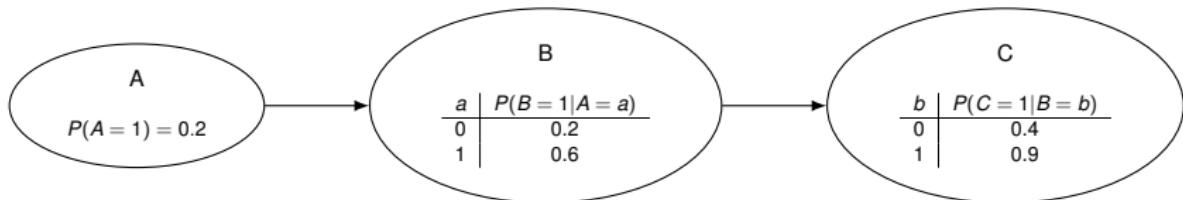


- ▶ $A \perp F ?$
- ▶ $A \perp F | D ?$
- ▶ $A \perp E ?$
- ▶ $A \perp E | D ?$
- ▶ $A \perp E | F ?$
- ▶ $A \perp E | B, D ?$

Réseaux bayésiens

Inférence

- ▶ Calcul d'une probabilité conditionnelle pour une certaine **évidence**
- ▶ Evidence = ensemble de variables à valeurs connues / observées

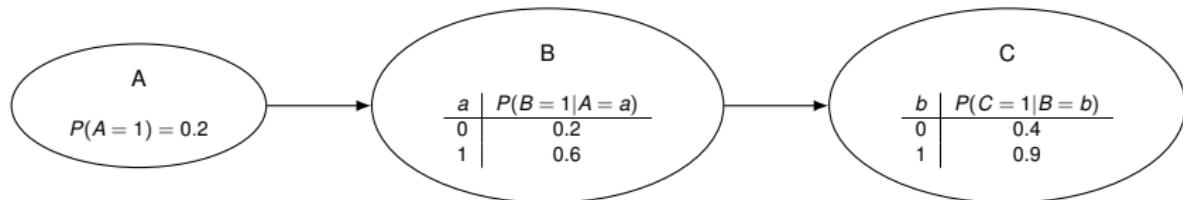


- ▶ $P(B = 1|A = 1)$?
- ▶ $P(C = 1|A = 1)$?

Réseaux bayésiens

Inférence

- ▶ Calcul d'une probabilité conditionnelle pour une certaine **évidence**
- ▶ Evidence = ensemble de variables à valeurs connues / observées

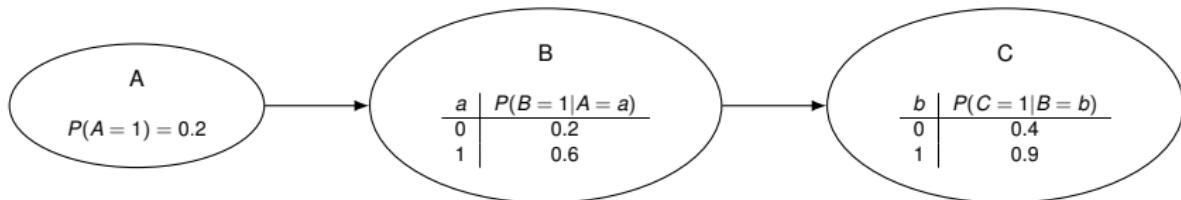


- ▶ $P(B = 1|A = 1) = 0.6 \rightarrow$ lecture directe depuis la table
- ▶ $P(C = 1|A = 1) ?$

Réseaux bayésiens

Inférence

- ▶ Calcul d'une probabilité conditionnelle pour une certaine **évidence**
- ▶ Evidence = ensemble de variables à valeurs connues / observées

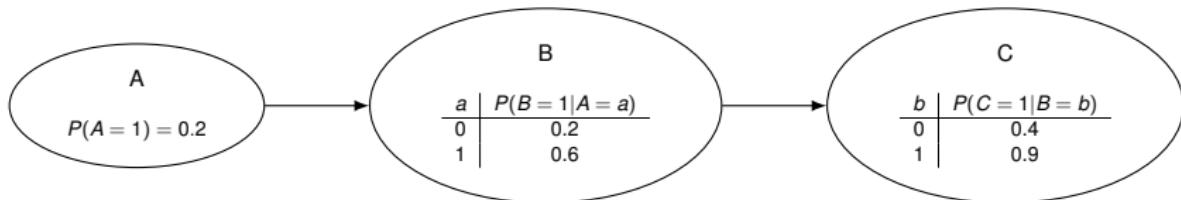


- ▶ $P(B = 1|A = 1) = 0.6 \rightarrow$ lecture directe depuis la table
- ▶ $P(C = 1|A = 1) \rightarrow$ utilisation des probabilités totales
 $= P(C = 1, B = 0|A = 1) + P(C = 1, B = 1|A = 1)$
 $= P(C = 1|A = 1, B = 0)P(B = 0|A = 1) + P(C = 1|A = 1, B = 1)P(B = 1|A = 1)$
 $= P(C = 1|B = 0)P(B = 0|A = 1) + P(C = 1|B = 1)P(B = 1|A = 1)$

Réseaux bayésiens

Inférence

- ▶ Calcul d'une probabilité conditionnelle pour une certaine **évidence**
- ▶ Evidence = ensemble de variables à valeurs connues / observées



- ▶ $P(B = 1|A = 1) = 0.6 \rightarrow$ lecture directe depuis la table
- ▶ $P(C = 1|A = 1) \rightarrow$ utilisation des probabilités totales
 $= P(C = 1, B = 0|A = 1) + P(C = 1, B = 1|A = 1)$
 $= P(C = 1|A = 1, B = 0)P(B = 0|A = 1) + P(C = 1|A = 1, B = 1)P(B = 1|A = 1)$
 $= P(C = 1|B = 0)P(B = 0|A = 1) + P(C = 1|B = 1)P(B = 1|A = 1)$

Règle de Bayes

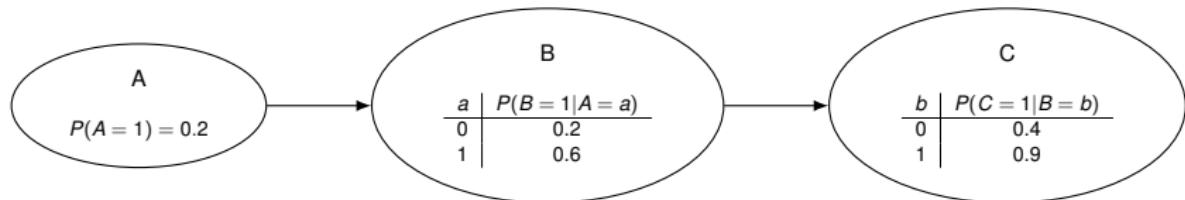
$$P(X = x | Y = y) = \frac{P(Y = y | X = x)P(X = x)}{P(Y = y)}$$

- ▶ $P(A = 1 | B = 1) ?$

Réseaux bayésiens

Inférence

- ▶ Calcul d'une probabilité conditionnelle pour une certaine **évidence**
- ▶ Evidence = ensemble de variables à valeurs connues / observées



- ▶ $P(B = 1|A = 1) = 0.6 \rightarrow$ lecture directe depuis la table
- ▶ $P(C = 1|A = 1) \rightarrow$ utilisation des probabilités totales
 $= P(C = 1, B = 0|A = 1) + P(C = 1, B = 1|A = 1)$
 $= P(C = 1|A = 1, B = 0)P(B = 0|A = 1) + P(C = 1|A = 1, B = 1)P(B = 1|A = 1)$
 $= P(C = 1|B = 0)P(B = 0|A = 1) + P(C = 1|B = 1)P(B = 1|A = 1)$

Règle de Bayes

$$P(X = x | Y = y) = \frac{P(Y = y | X = x)P(X = x)}{P(Y = y)}$$

- ▶ $P(A = 1 | B = 1)$
 $= \frac{0.6 \times 0.2}{P(B=1)} = \frac{0.12}{P(B=1|A=1)P(A=1)+P(B=1|A=0)P(A=0)} = \frac{0.12}{0.12+0.2 \times 0.8} = \frac{0.12}{0.28} \approx 0.43$

Réseaux bayésiens

Inférence en général

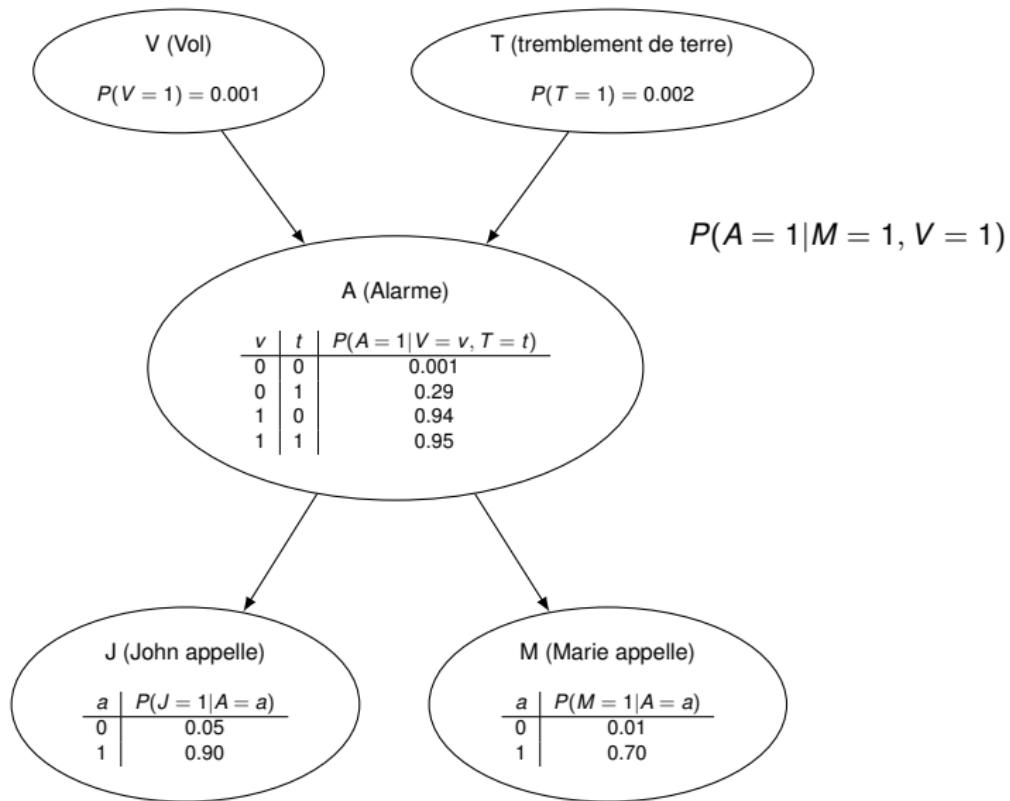
- ▶ Calcul d'une probabilité marginale : $P(X = 1) = \sum_{\mathbf{v}} P(X = 1, \mathbf{V} = \mathbf{v})$
- ▶ Calcul de $P(X = 1 | \mathbf{E} = \mathbf{e})$
- ▶ Si $\text{parents}(X) \subseteq \mathbf{E} \subseteq \text{predecesseurs}(X)$, alors lire la table des probabilités cond. de X
- ▶ Sinon, il y a d'autres variables \mathbf{V} et on passe par la loi jointe

$$\begin{aligned} P(X = 1 | \mathbf{E} = \mathbf{e}) &= \frac{P(X = 1, \mathbf{E} = \mathbf{e})}{P(\mathbf{E} = \mathbf{e})} = \sum_{\mathbf{v}} \frac{P(X = 1, \mathbf{V} = \mathbf{v}, \mathbf{E} = \mathbf{e})}{P(\mathbf{E} = \mathbf{e})} \\ &= \sum_{\mathbf{v}} P(X = 1, \mathbf{V} = \mathbf{v} | \mathbf{E} = \mathbf{e}) \\ &= \sum_{\mathbf{v}} P(X = 1 | \mathbf{V} = \mathbf{v}, \mathbf{E} = \mathbf{e}) P(\mathbf{V} = \mathbf{v} | \mathbf{E} = \mathbf{e}) \end{aligned}$$

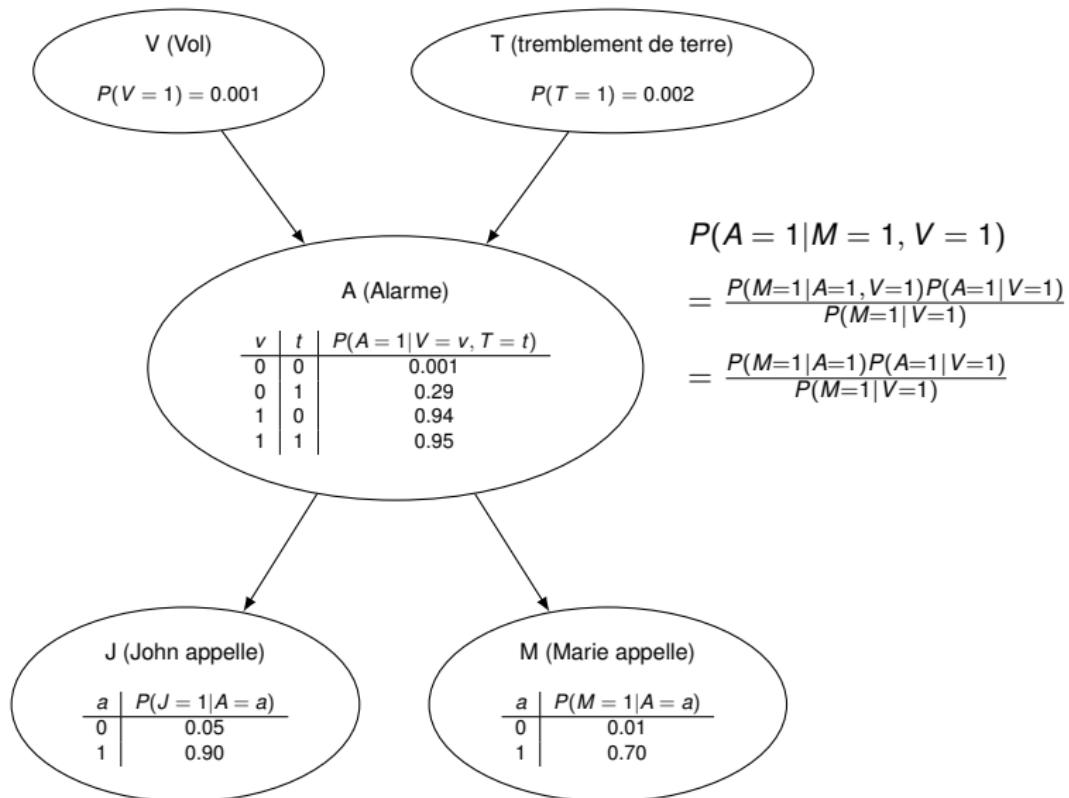
- ▶ La règle de Bayes permet d'inverser une condition :

$$P(X = x | (\mathbf{E}_1, \mathbf{E}_2) = (\mathbf{e}_1, \mathbf{e}_2)) = \frac{P(\mathbf{E}_1 = \mathbf{e}_1 | X = x, \mathbf{E}_2 = \mathbf{e}_2) P(X = x | \mathbf{E}_2 = \mathbf{e}_2)}{P(\mathbf{E}_1 = \mathbf{e}_1 | \mathbf{E}_2 = \mathbf{e}_2)}$$

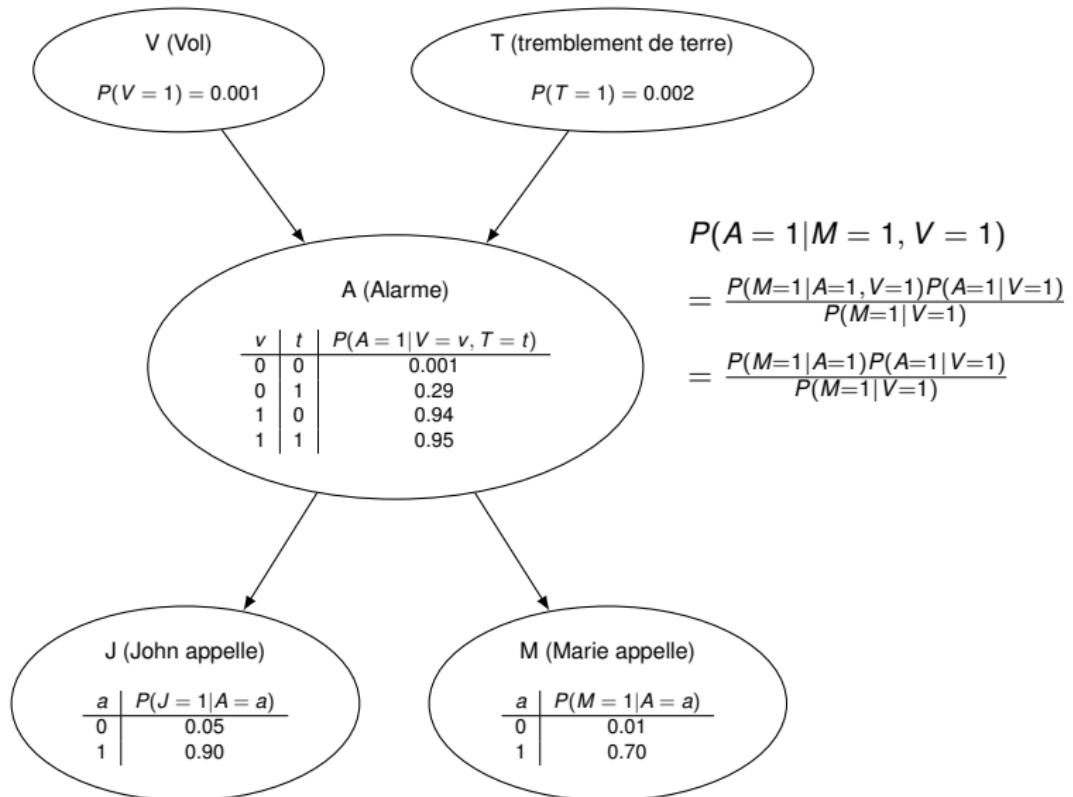
Réseaux bayésiens



Réseaux bayésiens



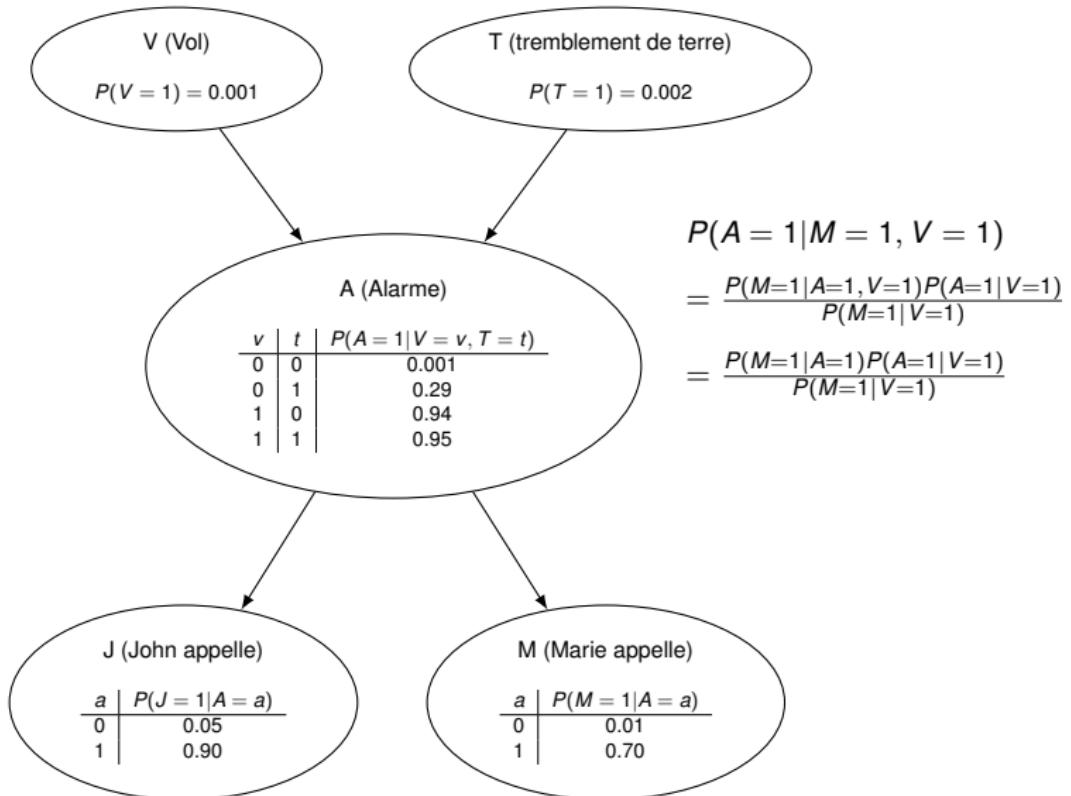
Réseaux bayésiens



$$\begin{aligned}
 & P(A = 1 | M = 1, V = 1) \\
 &= \frac{P(M=1|A=1, V=1)P(A=1|V=1)}{P(M=1|V=1)} \\
 &= \frac{P(M=1|A=1)P(A=1|V=1)}{P(M=1|V=1)}
 \end{aligned}$$

$$\begin{aligned}
 P(M = 1 | V = 1) &= P(M = 1 | V = 1, A = 1)P(A = 1 | V = 1) + P(M = 1 | V = 1, A = 0)P(A = 0 | V = 1) \\
 &= P(M = 1 | A = 1)P(A = 1 | V = 1) + P(M = 1 | A = 0)P(A = 0 | V = 1)
 \end{aligned}$$

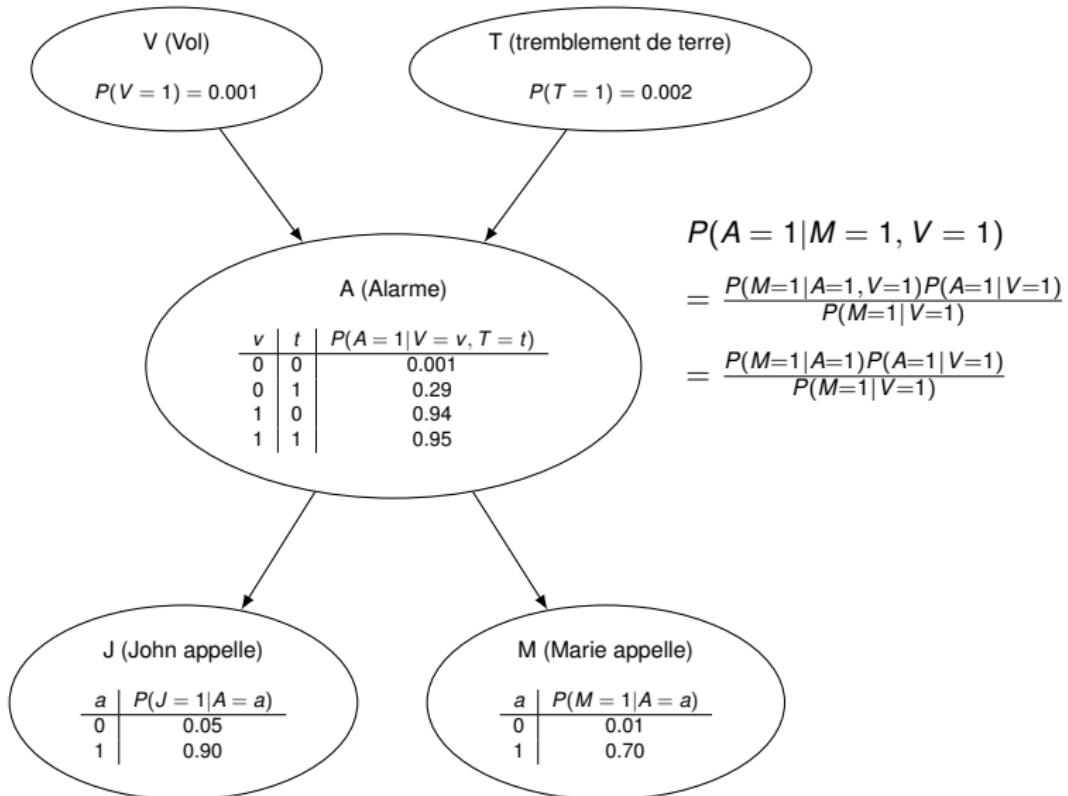
Réseaux bayésiens



$$\begin{aligned}
 P(A = 1 | M = 1, V = 1) \\
 &= \frac{P(M=1|A=1, V=1)P(A=1|V=1)}{P(M=1|V=1)} \\
 &= \frac{P(M=1|A=1)P(A=1|V=1)}{P(M=1|V=1)}
 \end{aligned}$$

$$\begin{aligned}
 P(M = 1 | V = 1) &= P(M = 1 | V = 1, A = 1)P(A = 1 | V = 1) + P(M = 1 | V = 1, A = 0)P(A = 0 | V = 1) \\
 &= P(M = 1 | A = 1)P(A = 1 | V = 1) + P(M = 1 | A = 0)P(A = 0 | V = 1) \\
 P(A = 1 | V = 1) &= P(A = 1 | V = 1, T = 0)P(T = 0) + P(A = 1 | V = 1, T = 1)P(T = 1)
 \end{aligned}$$

Réseaux bayésiens

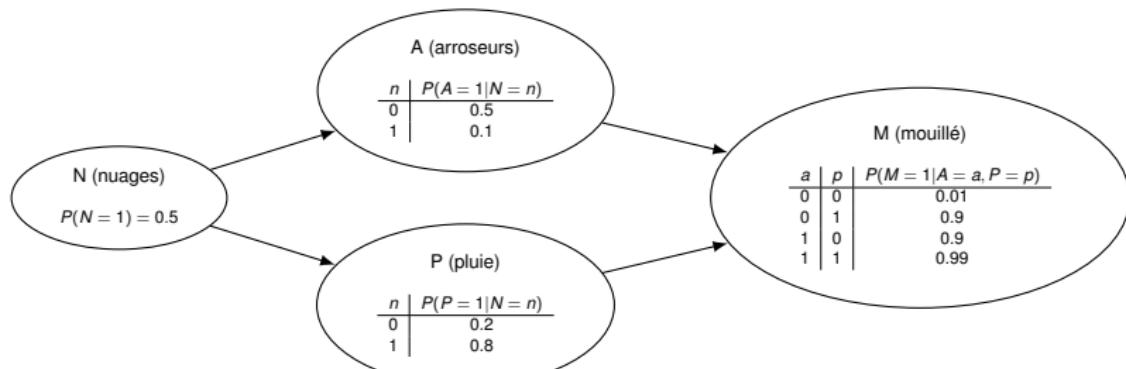


$$\begin{aligned}
 P(M = 1 | V = 1) &= P(M = 1 | V = 1, A = 1)P(A = 1 | V = 1) + P(M = 1 | V = 1, A = 0)P(A = 0 | V = 1) \\
 &= P(M = 1 | A = 1)P(A = 1 | V = 1) + P(M = 1 | A = 0)P(A = 0 | V = 1) \\
 P(A = 1 | V = 1) &= P(A = 1 | V = 1, T = 0)P(T = 0) + P(A = 1 | V = 1, T = 1)P(T = 1) \\
 P(A = 0 | V = 1) &= 1 - P(A = 1 | V = 1)
 \end{aligned}$$

Réseaux bayésiens

Inférence approximative

- ▶ Inférence exacte souvent trop complexe (NP-difficile)
- ▶ Approximation par échantillonnage (en commençant par les nœuds dont les parents sont inexistant ou déjà échantillonnés)
- ▶ $P(\mathbf{X} = \mathbf{x}) = \mathbb{E}[\mathbb{I}(\mathbf{X} = \mathbf{x})] \approx \frac{\text{nb d'échantillons avec } \mathbf{X}=\mathbf{x}}{\text{nb total d'échantillons}}$

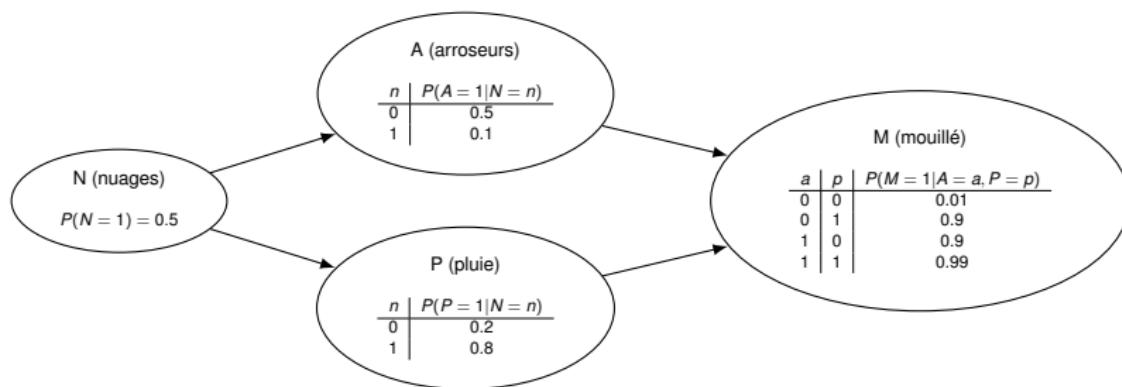


$$P(N = 1, A = 0, P = 1, M = 1) ?$$

Réseaux bayésiens

Inférence approximative de lois conditionnelles

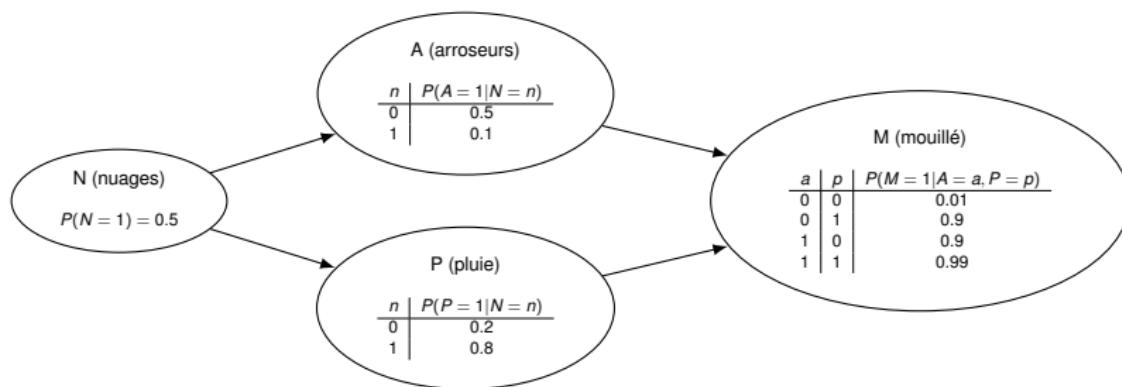
- $P(X = x | \mathbf{E} = \mathbf{e}) \approx \frac{\text{nb d'échantillons avec } X=x \text{ et } \mathbf{E}=\mathbf{e}}{\text{nb d'échantillons avec } \mathbf{E}=\mathbf{e}}$
- $P(M = 1 | N = 0)$?



Réseaux bayésiens

Inférence approximative de lois conditionnelles

- ▶ $P(X = x | \mathbf{E} = \mathbf{e}) \approx \frac{\text{nb d'échantillons avec } X=x \text{ et } \mathbf{E}=\mathbf{e}}{\text{nb d'échantillons avec } \mathbf{E}=\mathbf{e}}$
- ▶ On rejette certains échantillons et on perd beaucoup de temps (surtout si $\mathbf{E} = \mathbf{e}$ est rare)
 $P(P = 1 | A = 1, M = 1) ?$

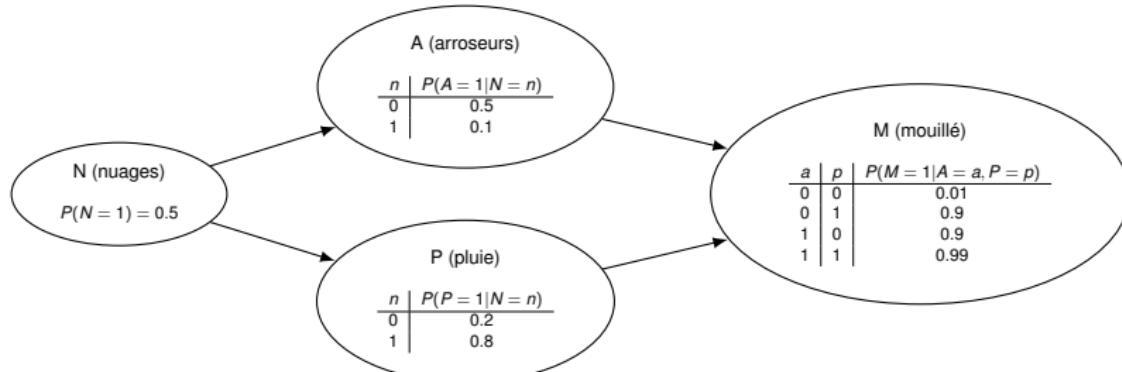


Réseaux bayésiens

Inférence approximative de lois conditionnelles

- ▶ $P(X = x | \mathbf{E} = \mathbf{e}) \approx \frac{\text{nb d'échantillons avec } X=x \text{ et } \mathbf{E}=\mathbf{e}}{\text{nb d'échantillons avec } \mathbf{E}=\mathbf{e}}$
- ▶ On rejette certains échantillons et on perd beaucoup de temps (surtout si $\mathbf{E} = \mathbf{e}$ est rare)
- ▶ Pondération par la vraisemblance : forcer l'évidence $\mathbf{E} = \mathbf{e}$ au lieu de tirer aléatoirement \mathbf{E} et pondérer les échantillons par la vraisemblance v_i d'avoir $\mathbf{E} = \mathbf{e}$ pour chaque échantillon i :

$$P(X = x | \mathbf{E} = \mathbf{e}) \approx \frac{\sum_i v_i \mathbb{I}(X = x)}{\sum_i v_i}$$



Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

Classification multi-classe

Méthodes non linéaires

Méthodes à noyaux

Apprentissage automatique

Définition de A. Samuel (1959)

Apprentissage automatique (*Machine Learning*) : champ d'étude donnant aux ordinateurs la capacité d'apprendre sans être explicitement programmés

Apprentissage automatique

Définition de A. Samuel (1959)

Apprentissage automatique (*Machine Learning*) : champ d'étude donnant aux ordinateurs la capacité d'apprendre sans être explicitement programmés

Apprentissage par l'exemple

- ▶ Méthode privilégiée par l'apprentissage automatique
- ▶ Exemple : reconnaissance de chiffres manuscrits

Apprentissage automatique

Définition de A. Samuel (1959)

Apprentissage automatique (*Machine Learning*) : champ d'étude donnant aux ordinateurs la capacité d'apprendre sans être explicitement programmés

Apprentissage par l'exemple

- ▶ Méthode privilégiée par l'apprentissage automatique
- ▶ Exemple : reconnaissance de chiffres manuscrits

Définition

Apprendre = Acquérir la capacité de généraliser à partir d'exemples

Généraliser = Prédire correctement un phénomène non observé auparavant

Apprentissage automatique

Définition de A. Samuel (1959)

Apprentissage automatique (*Machine Learning*) : champ d'étude donnant aux ordinateurs la capacité d'apprendre sans être explicitement programmés

Apprentissage par l'exemple

- ▶ Méthode privilégiée par l'apprentissage automatique
- ▶ Exemple : reconnaissance de chiffres manuscrits

Définition

Apprendre = Acquérir la capacité de généraliser à partir d'exemples

Généraliser = Prédire correctement un phénomène non observé auparavant

Apprendre par cœur ne permet pas de généraliser

- ▶ Exemple : apprendre les préférences d'un utilisateur pour les touches du clavier

Motivations / Applications

Vision par ordinateur

- ▶ Reconnaissance de formes (ex. codes postaux)
- ▶ Suivi de trajectoires
- ▶ Segmentation d'images / de vidéo

Bioinformatique

- ▶ Prédiction de la structure secondaire des protéines
- ▶ Philogénie
- ▶ ...

Biomédical

- ▶ Aide au diagnostic, pronostics
- ▶ Imagerie médicale

Physique / ingénierie / automatique

- ▶ Modélisation de processus complexes
- ▶ Simulation temps réel
- ▶ Lois de commandes

Informatique

- ▶ Détection d'intrusion
- ▶ Détection de SPAM
- ▶ Applications web (propositions orthographiques, détection de langue, publicité ciblée)

Robotique

- ▶ Perception artificielle
- ▶ Prise de décisions

Fouille de données

- ▶ Extraction de connaissances

Séries temporelles

- ▶ Finance
- ▶ Météo
- ▶ Charge du réseau (EDF)

Google Prediction API



Google Prediction API

Google's cloud-based machine learning tools can help analyze your data to add the following features to your applications:



Customer sentiment analysis



Spam detection



Message routing decisions



Upsell opportunity analysis



Document and email classification



Diagnostics



Churn analysis



Suspicious activity identification



Recommendation systems



And much more...

Microsoft Azure Machine Learning API



Recommandations

Aidez les clients à découvrir des articles de votre catalogue et améliorez les ventes dans vos boutiques en ligne.

[Inscrivez-vous dans le portail Microsoft Azure Classic](#)



Analyse de texte

Effectuez l'analyse des sentiments et l'extraction des phrases clés sur du texte non structuré (en anglais uniquement).

[Inscrivez-vous dans le portail Microsoft Azure Classic](#)



Prédiction de la désinscription des clients

Prévoyez la probabilité qu'une relation se termine entre un client et une société/un service.

[Inscrivez-vous dans le portail Microsoft Azure Classic](#)



API Face

Conçu pour détecter et reconnaître les visages dans les images.

[Inscrivez-vous dans le portail](#)



API Speech

Ajoutez des actions vocales à vos applications.

[Inscrivez-vous dans le portail](#)



API Vision

Conçues pour le contenu visuel, tel que l'analyse des images, la génération de miniatures et l'extraction de caractères.

[Inscrivez-vous dans le portail](#)

Représentation des exemples

Un exemple est un couple $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$

Espace de représentation : \mathcal{X}

- ▶ Les objets pour lesquels on souhaite prédire quelque chose sont représentés par un ensemble de descripteurs (*features*) ou attributs

$$\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathcal{X} \subseteq \mathbb{R}^d$$

- ▶ d : dimension de l'espace de représentation

Ensemble des étiquettes possibles : \mathcal{Y}

Etant donné la représentation \mathbf{x} d'un objet, on cherche à prédire son **étiquette** $y \in \mathcal{Y}$

- ▶ Discrimination (ou classification) binaire : $\mathcal{Y} = \{-1, +1\}$
- ▶ Discrimination multi-classe : $\mathcal{Y} = \{1, \dots, Q\}$
- ▶ Régression : $\mathcal{Y} = \mathbb{R}$

Le choix des descripteurs est pertinent si ces variables sont liées à l'étiquette

De quoi dispose-t-on pour prédire ?

Ensemble des exemples disponibles S

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$$

- ▶ S est la **base d'apprentissage**
- ▶ m : taille de la base d'apprentissage = nombre d'exemples disponibles
- ▶ S contient (souvent) les seules informations disponibles

Modèle de prédiction f

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

- ▶ Prédire = calculer la sortie $f(\mathbf{x})$ du modèle
- ▶ \mathcal{F} : ensemble des modèles possibles ($f \in \mathcal{F}$)
 - ▶ \mathcal{F} est une classe de fonctions (= famille de fonctions = ensemble de fonctions)
 - ▶ La recherche du modèle de prédiction f est limitée à l'ensemble \mathcal{F}

Comment apprendre à prédire ?

Algorithme d'apprentissage \mathcal{A}

$$\mathcal{A} : (\mathcal{F}, S) \mapsto \hat{f} \in \mathcal{F}$$

Apprendre = sélectionner un modèle \hat{f} parmi les modèles f de \mathcal{F} après observation d'un ensemble d'exemples S

- ▶ But : sélectionner un modèle de prédiction \hat{f} capable de "prédire correctement" l'étiquette y pour tout $\mathbf{x} \in \mathcal{X}$ et pas uniquement pour les exemples de la base d'apprentissage

Comment apprendre à prédire ?

Algorithme d'apprentissage \mathcal{A}

$$\mathcal{A} : (\mathcal{F}, S) \mapsto \hat{f} \in \mathcal{F}$$

Apprendre = sélectionner un modèle \hat{f} parmi les modèles f de \mathcal{F} après observation d'un ensemble d'exemples S

- ▶ But : sélectionner un modèle de prédiction \hat{f} capable de "prédire correctement" l'étiquette y pour tout $\mathbf{x} \in \mathcal{X}$ et pas uniquement pour les exemples de la base d'apprentissage

Généralisation

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{f}(\mathbf{x}) \stackrel{?}{=} y$$

- ▶ Généralisation : capacité de \hat{f} à donner des prédictions $\hat{f}(\mathbf{x})$ proches de y pour $\mathbf{x} \notin S$

Paramètres et hyperparamètres

Paramètres

- ▶ Ensemble des variables paramétrant les fonctions f d'une classe \mathcal{F} donnée
- ▶ Exemples :

$$\mathcal{F}_1 = \left\{ f \mid f(x) = x^\theta, \theta \in \{1, 2, 3\} \right\}$$

$$\mathcal{F}_2 = \left\{ f \mid f(x) = \theta_1 x + \theta_2, (\theta_1, \theta_2) \in \mathbb{R}^2 \right\}$$

Apprendre = déterminer les valeurs des paramètres θ

Paramètres et hyperparamètres

Paramètres

- ▶ Ensemble des variables paramétrant les fonctions f d'une classe \mathcal{F} donnée
- ▶ Exemples :

$$\mathcal{F}_1 = \left\{ f \mid f(x) = x^\theta, \theta \in \{1, 2, 3\} \right\}$$

$$\mathcal{F}_2 = \left\{ f \mid f(x) = \theta_1 x + \theta_2, (\theta_1, \theta_2) \in \mathbb{R}^2 \right\}$$

Apprendre = déterminer les valeurs des paramètres θ

Hyperparamètres

- ▶ Ensemble des variables paramétrant un algorithme d'apprentissage ou une classe de fonctions \mathcal{F}
- ▶ Exemple : \mathcal{F}^γ est l'ensembles des fonctions linéaires de x^γ

$$\mathcal{F}^\gamma = \{f \mid f(x) = \theta x^\gamma, \theta \in \mathbb{R}\}$$

- ▶ Autre exemple : l'algorithme \mathcal{A}^γ stoppe après γ itérations
- ▶ L'apprentissage des paramètres θ peut être réalisé pour différentes valeurs de γ
- ▶ Chaque valeur de γ conduit à un apprentissage et un modèle \hat{f} différents

Comment valider l'apprentissage ?

Validation

- ▶ Etape destinée au **réglage des hyperparamètres** γ (sélection de modèle)
- ▶ Evaluation de la qualité d'une valeur de γ *après* l'apprentissage du modèle \hat{f}
- ▶ En pratique : évaluer la qualité du modèle \hat{f} sur des données supplémentaires dites de validation

Comment valider l'apprentissage ?

Validation

- ▶ Etape destinée au **réglage des hyperparamètres γ** (sélection de modèle)
- ▶ Evaluation de la qualité d'une valeur de γ *après* l'apprentissage du modèle \hat{f}
- ▶ En pratique : évaluer la qualité du modèle \hat{f} sur des données supplémentaires dites de validation

Test

- ▶ **Etape finale** de la procédure destinée à **estimer la qualité du modèle \hat{f}** *après* réglage des éventuels hyperparamètres
- ▶ La qualité du modèle correspond à sa **capacité à généraliser**
- ▶ Cette estimation est indispensable pour déterminer si le modèle est utilisable et le risque lié à cette utilisation
- ▶ En pratique (dans les cas simples) : estimation de la qualité du modèle \hat{f} sur des données supplémentaires dites de test

Comment valider l'apprentissage ?

Validation

- ▶ Etape destinée au **réglage des hyperparamètres γ** (sélection de modèle)
- ▶ Evaluation de la qualité d'une valeur de γ *après* l'apprentissage du modèle \hat{f}
- ▶ En pratique : évaluer la qualité du modèle \hat{f} sur des données supplémentaires dites de validation

Test

- ▶ **Etape finale** de la procédure destinée à **estimer la qualité du modèle \hat{f}** *après* réglage des éventuels hyperparamètres
- ▶ La qualité du modèle correspond à sa **capacité à généraliser**
- ▶ Cette estimation est indispensable pour déterminer si le modèle est utilisable et le risque lié à cette utilisation
- ▶ En pratique (dans les cas simples) : estimation de la qualité du modèle \hat{f} sur des données supplémentaires dites de test

NE PAS utiliser les données de test pour régler les hyperparamètres
(cela conduirait à une surestimation des performances du modèle)

Comment valider l'apprentissage ?

Validation croisée

- ▶ Techniques d'estimation de la qualité d'un modèle (en test ou validation)
- ▶ Pour les cas où le nombre de données est trop faible

Comment valider l'apprentissage ?

Validation croisée

- ▶ Techniques d'estimation de la qualité d'un modèle (en test ou validation)
- ▶ Pour les cas où le nombre de données est trop faible

Validation croisée "K-fold"

- ▶ Découper les données en K sous-ensembles (typiquement $K = 5$ ou $K = 7$)
- ▶ Réaliser K apprentissages sur $K - 1$ sous-ensembles en mesurant la qualité de chaque modèle sur le sous-ensemble restant
- ▶ La moyenne des K mesures de qualité donne une estimation de la qualité utilisant toutes les données

Validation croisée "Leave-one-out (LOO)"

- ▶ Cas limite de la méthode K-fold avec K égal au nombre de données
- ▶ Autant d'apprentissages que de données, la qualité de chaque apprentissage est mesurée sur un seul exemple
- ▶ Très coûteux en temps de calcul

Modélisation probabiliste de l'apprentissage

Cadre probabiliste pour l'apprentissage

- ▶ L'espace $(\mathcal{X}, \mathcal{Y})$ est muni d'une mesure de probabilité P **inconnue**
- ▶ $(X, Y) \in (\mathcal{X}, \mathcal{Y})$: couple de variables aléatoires de loi P ($= P_{X,Y}$)
- ▶ Jeu de données $S = \{(x_i, y_i)\}_{1 \leq i \leq m}$:
 m réalisations **indépendantes et identiquement distribuées** (i.i.d.) de (X, Y)
- ▶ Notation : X et x_i peuvent être des vecteurs de dimension d

Modélisation probabiliste de l'apprentissage

Cadre probabiliste pour l'apprentissage

- ▶ L'espace $(\mathcal{X}, \mathcal{Y})$ est muni d'une mesure de probabilité P **inconnue**
- ▶ $(X, Y) \in (\mathcal{X}, \mathcal{Y})$: couple de variables aléatoires de loi P ($= P_{X,Y}$)
- ▶ Jeu de données $S = \{(x_i, y_i)\}_{1 \leq i \leq m}$:
 m réalisations **indépendantes et identiquement distribuées** (i.i.d.) de (X, Y)
- ▶ Notation : X et x_i peuvent être des vecteurs de dimension d

Remarques

- ▶ **indépendantes** : chaque exemple (x_i, y_i) apporte le maximum d'information
- ▶ **identiquement distribuées** : chaque exemple (x_i, y_i) suit la loi P inconnue et toute observation future du couple (X, Y) est supposée suivre la même loi
la base d'apprentissage est donc sensée être représentative du problème
- ▶ **Seules les régularités des données peuvent être apprises**
- ▶ **On ne peut pas prédire l'imprévisible**
(ex : crack boursier)

Modélisation probabiliste de l'apprentissage

2 cas à distinguer

Cas déterministe

- ▶ Il existe une fonction t telle que $Y = t(X)$
- ▶ $\forall x, P(Y = t(X) | X = x) = 1$
- ▶ Y est entièrement déterminé par X
(X contient suffisamment d'information pour déterminer Y sans ambiguïté)
- ▶ Exemple :
$$Y = +1 \text{ si } X \geq 5, -1 \text{ sinon}$$

Modélisation probabiliste de l'apprentissage

2 cas à distinguer

Cas déterministe

- ▶ Il existe une fonction t telle que $Y = t(X)$
- ▶ $\forall x, P(Y = t(X) | X = x) = 1$
- ▶ Y est entièrement déterminé par X
(X contient suffisamment d'information pour déterminer Y sans ambiguïté)
- ▶ Exemple :
$$Y = +1 \text{ si } X \geq 5, -1 \text{ sinon}$$

Cas stochastique

- ▶ $\forall t, P(Y = t(X)) \neq 1$
 $\Rightarrow P(Y = t(X) | X = x) \neq 1$ pour un ensemble de x
- ▶ Pour $X = x$ donné, il existe plusieurs Y plausibles (+ ou - probables)
- ▶ Interprétations de l'aléatoire : bruit de mesure, manque d'information, simplification de modèle
Exemple : classification homme/femme par rapport à la taille

Mesures de performance et apprentissage

Apprentissage d'un modèle de prédiction

- ▶ **Modèle de prédiction** : $f : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ **Algorithmе d'apprentissage** : $\mathcal{A} : (\mathcal{F}, S) \mapsto \hat{f} \in \mathcal{F}$

Mesures de performance et apprentissage

Apprentissage d'un modèle de prédiction

- ▶ **Modèle de prédiction** : $f : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ **Algorithmе d'apprentissage** : $\mathcal{A} : (\mathcal{F}, S) \mapsto \hat{f} \in \mathcal{F}$

Mesure de l'erreur de prédiction sur un exemple : la fonction de perte

- ▶ **Fonction de perte** $\ell : (\mathcal{F}, \mathcal{X}, \mathcal{Y}) \rightarrow \mathbb{R}^+$
- ▶ Propriété : $f(x) = y \Rightarrow \ell(f, x, y) = 0$

Mesures de performance et apprentissage

Apprentissage d'un modèle de prédiction

- ▶ **Modèle de prédiction** : $f : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ **Algorithme d'apprentissage** : $\mathcal{A} : (\mathcal{F}, S) \mapsto \hat{f} \in \mathcal{F}$

Mesure de l'erreur de prédiction sur un exemple : la fonction de perte

- ▶ **Fonction de perte** $\ell : (\mathcal{F}, \mathcal{X}, \mathcal{Y}) \rightarrow \mathbb{R}^+$
- ▶ Propriété : $f(x) = y \Rightarrow \ell(f, x, y) = 0$

Le risque d'un modèle de prédiction

Risque de f = **erreur de généralisation** de f = espérance de la fonction de perte calculée en f :

$$R(f) = \mathbb{E}_{(X, Y)} [\ell(f, X, Y)]$$

Mesures de performance et apprentissage

Apprentissage d'un modèle de prédiction

- ▶ **Modèle de prédiction** : $f : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ **Algorithme d'apprentissage** : $\mathcal{A} : (\mathcal{F}, S) \mapsto \hat{f} \in \mathcal{F}$

Mesure de l'erreur de prédiction sur un exemple : la fonction de perte

- ▶ **Fonction de perte** $\ell : (\mathcal{F}, \mathcal{X}, \mathcal{Y}) \rightarrow \mathbb{R}^+$
- ▶ Propriété : $f(x) = y \Rightarrow \ell(f, x, y) = 0$

Le risque d'un modèle de prédiction

Risque de f = **erreur de généralisation** de f = espérance de la fonction de perte calculée en f :

$$R(f) = \mathbb{E}_{(X, Y)} [\ell(f, X, Y)]$$

Risque = erreur (mesurée par la fonction de perte) probable ou "en moyenne" sur un exemple tiré aléatoirement

But de l'apprentissage : minimiser le risque

Risque de f = erreur de généralisation de f :

$$R(f) = \mathbb{E}_{(X, Y)} [\ell(f, X, Y)]$$

But de l'apprentissage : minimiser le risque

Risque de f = erreur de généralisation de f :

$$R(f) = \mathbb{E}_{(X, Y)} [\ell(f, X, Y)]$$

Minimisation du risque

- ▶ Minimiser le risque revient à minimiser l'erreur probable de prédiction
- ▶ Apprentissage idéal : retenir le modèle f avec le risque $R(f)$ minimal
- ▶ Mais le risque $R(f)$ dépend de la loi P du couple (X, Y) qui est inconnue
- ▶ Approche : minimiser une estimation du risque

But de l'apprentissage : minimiser le risque

Risque de f = erreur de généralisation de f :

$$R(f) = \mathbb{E}_{(X, Y)} [\ell(f, X, Y)]$$

Minimisation du risque

- ▶ Minimiser le risque revient à minimiser l'erreur probable de prédiction
- ▶ Apprentissage idéal : retenir le modèle f avec le risque $R(f)$ minimal
- ▶ Mais le risque $R(f)$ dépend de la loi P du couple (X, Y) qui est inconnue
- ▶ Approche : minimiser une estimation du risque

Estimer le risque à partir de données empiriques

- ▶ **Erreur d'apprentissage** (= erreur empirique = risque empirique) :

$$R_{emp}(f) = \frac{1}{m} \sum_{i=1}^m \ell(f, x_i, y_i)$$

- ▶ Moyenne empirique de la fonction de perte sur la base d'apprentissage

Principe de minimisation du risque empirique

Principe inductif ERM (*empirical risk minimization*)

- ▶ Choisir le modèle qui minimise $R_{emp}(f)$ au lieu de $R(f)$
- ▶ Soit l'algorithme :

$$\hat{f} = \arg \min_{f \in \mathcal{F}} R_{emp}(f)$$

\hat{f} : modèle de prédiction retenu par l'algorithme d'apprentissage parmi l'ensemble des modèles possibles \mathcal{F}

Remarques

- ▶ Les algorithmes d'apprentissage prennent souvent la forme de problèmes de minimisation
- ▶ On appelle *algorithme d'apprentissage* soit la définition du problème d'optimisation (ci-dessus) soit l'algorithme permettant de le résoudre
- ▶ Le principe ERM est simple mais malgré tout souvent efficace (si \mathcal{F} est choisi correctement !)

Discrimination binaire : $\mathcal{Y} = \{-1, +1\}$

- ▶ Prédire Y revient à classer X ; f est appelé *classifieur*

Risque de classification

- ▶ Fonction indicatrice : $\mathbb{I}(A) = 1$ si A est observé, 0 sinon
- ▶ Fonction de perte standard : $\ell(f, X, Y) = \mathbb{I}(f(X) \neq Y)$
- ▶ **Risque = probabilité de mal classer un exemple** = "taux d'erreur" :

$$R(f) = \mathbb{E}_{(X, Y)}[\mathbb{I}(f(X) \neq Y)] = P(f(X) \neq Y)$$

- ▶ Taux de reconnaissance = $1 - P(f(X) \neq Y)$
- ▶ Sur la base d'apprentissage :
taux de reconnaissance = $1 - R_{emp}(f) = 1 - \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) \neq y_i)$

Discrimination binaire : $\mathcal{Y} = \{-1, +1\}$

- ▶ Prédire Y revient à classer X ; f est appelé *classifieur*

Risque de classification

- ▶ Fonction indicatrice : $\mathbb{I}(A) = 1$ si A est observé, 0 sinon
- ▶ Fonction de perte standard : $\ell(f, X, Y) = \mathbb{I}(f(X) \neq Y)$
- ▶ **Risque = probabilité de mal classer un exemple** = "taux d'erreur" :

$$R(f) = \mathbb{E}_{(X, Y)}[\mathbb{I}(f(X) \neq Y)] = P(f(X) \neq Y)$$

- ▶ Taux de reconnaissance = $1 - R(f)$
- ▶ Sur la base d'apprentissage :
taux de reconnaissance = $1 - R_{emp}(f) = 1 - \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) \neq y_i)$

Classifieur optimal : le classifieur de Bayes

Associer l'exemple à la classe la plus probable :

$$f_{Bayes}(x) = \begin{cases} +1, & \text{si } P(Y = 1|X = x) \geq P(Y = -1|X = x) \\ -1, & \text{si } P(Y = 1|X = x) < P(Y = -1|X = x) \end{cases}$$

Non applicable en pratique sans la connaissance de P

Discrimination binaire : $\mathcal{Y} = \{-1, +1\}$

Classifieur optimal : le classifieur de Bayes

Associer l'exemple à la classe la plus probable :

$$f_{Bayes}(x) = \begin{cases} +1, & \text{si } P(Y = 1|X = x) \geq P(Y = -1|X = x) \\ -1, & \text{si } P(Y = 1|X = x) < P(Y = -1|X = x) \end{cases}$$

Discrimination binaire : $\mathcal{Y} = \{-1, +1\}$

Classifieur optimal : le classifieur de Bayes

Associer l'exemple à la classe la plus probable :

$$f_{Bayes}(x) = \begin{cases} +1, & \text{si } P(Y = 1|X = x) \geq P(Y = -1|X = x) \\ -1, & \text{si } P(Y = 1|X = x) < P(Y = -1|X = x) \end{cases}$$

Risque de Bayes

$$R^* = R(f_{Bayes}) = \inf_f R(f)$$

- ▶ Risque minimal sur l'ensemble de toutes les fonctions $f : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ Cas déterministe : $P(Y = 1|X = x) \in \{0, 1\} \Rightarrow R^* = 0$
(pas d'erreur possible, le plus probable a toujours une probabilité de 1)
- ▶ Cas stochastique : $P(Y = 1|X = x) \in [0, 1] \Rightarrow R^* > 0$
(quelque soit la réponse, il existe une probabilité non nulle de se tromper)

Discrimination binaire : $\mathcal{Y} = \{-1, +1\}$

Classifieur optimal : le classifieur de Bayes

- Le classifieur optimal f^* est celui qui minimise le risque : $f^* = \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} R(f)$
- Classifieur de Bayes : $\forall x \in \mathcal{X}, f_{Bayes}(x) = 1 \Leftrightarrow P(Y = 1 | X = x) \geq \frac{1}{2}$

Preuve de $f^* = f_{Bayes}$

- Pour tout $x \in \mathcal{X}$, et n'importe quel classifieur $f : \mathcal{X} \rightarrow \mathcal{Y}$:

$$\begin{aligned} P(f(X) \neq Y | X = x) &= P(f(X) = 1, Y = -1 | X = x) + P(f(X) = -1, Y = 1 | X = x) \\ &= \mathbb{I}(f(x) = 1)P(Y = -1 | X = x) + \mathbb{I}(f(x) = -1)P(Y = 1 | X = x) \\ &= \mathbb{I}(f(x) = 1)[1 - P(Y = 1 | X = x)] + [1 - \mathbb{I}(f(x) = 1)]P(Y = 1 | X = x) \\ &= [1 - 2P(Y = 1 | X = x)]\mathbb{I}(f(x) = 1) + P(Y = 1 | X = x) \end{aligned}$$

- Donc

$$\begin{aligned} \Delta &= P(f(X) \neq Y | X = x) - P(f_{Bayes}(X) \neq Y | X = x) \\ &= \underbrace{[1 - 2P(Y = 1 | X = x)]}_{a} \underbrace{[\mathbb{I}(f(x) = 1) - \mathbb{I}(f_{Bayes}(x) = 1)]}_{b} \end{aligned}$$

- $\Delta = ab \geq 0$ pour tout f :

Cas 1 : $a \geq 0 \Rightarrow P(Y = 1 | X = x) \leq \frac{1}{2} \Rightarrow f_{Bayes}(x) = -1$
 $\Rightarrow b = \mathbb{I}(f(x) = 1) \geq 0 \Rightarrow \Delta \geq 0$

Cas 2 : $a \leq 0 \Rightarrow P(Y = 1 | X = x) \geq \frac{1}{2} \Rightarrow f_{Bayes}(x) = 1$
 $\Rightarrow b = \mathbb{I}(f(x) = 1) - 1 \leq 0 \Rightarrow \delta \geq 0$

- Donc $R(f) = \mathbb{E}_X P(f(X) \neq Y | X = x) \geq \mathbb{E}_X P(f_{Bayes}(X) \neq Y | X = x) = R(f_{Bayes})$

Régression : $\mathcal{Y} \subseteq \mathbb{R}$

- ▶ $X \in \mathbb{R}^d$: vecteur aléatoire continu (entrées, régresseurs)
- ▶ Y : v.a. continue (sortie)
- ▶ Loi jointe P à densité $p(x, y) = p(y|x)p(x)$

Risque de régression

- ▶ Fonction de perte quadratique : $\ell(f, X, Y) = (Y - f(X))^2$
- ▶ **Risque = erreur quadratique moyenne** (EQM ou MSE en anglais) :

$$R(f) = \mathbb{E}_{(X, Y)}[(Y - f(X))^2] = \int_{\mathcal{X} \times \mathcal{Y}} (y - f(x))^2 p(x, y) dx dy$$

Modèle optimal : la fonction de régression

La fonction de régression est l'espérance conditionnelle de Y sachant X :

$$\forall x \in \mathcal{X}, \quad f_{reg}(x) = \mathbb{E}_{Y|X=x}[Y|X = x]$$

Régression : $\mathcal{Y} \subseteq \mathbb{R}$

Modèle optimal = la fonction de régression

- Le modèle optimal f^* est celui qui minimise le risque :

$$f^* = \arg \min_{f: \mathcal{X} \rightarrow \mathbb{R}} R(f)$$

- Fonction de régression : $\forall x \in \mathcal{X}, f_{reg}(x) = \mathbb{E}_{Y|X=x}[Y|X=x]$

Preuve de $f^* = f_{reg}$

- $R(f) = \mathbb{E}_{(X,Y)}[(Y - f(X))^2] = \mathbb{E}_X \mathbb{E}_{Y|X=x}[(Y - f(X))^2 | X = x]$
- Minimiser $R(f)$ revient à trouver, pour tout x , la valeur $y = f^*(x)$ qui minimise $\mathbb{E}_{Y|X=x}[(Y - f(X))^2 | X = x]$:

$$\forall x \in \mathcal{X}, f^*(x) = \arg \min_{y \in \mathbb{R}} \mathbb{E}_{Y|X=x}[(Y - y)^2 | X = x]$$

$$\Leftrightarrow \forall x \in \mathcal{X}, f^*(x) = \arg \min_{y \in \mathbb{R}} \mathbb{E}_{Y|X=x}[Y^2 | X = x] - \underbrace{2y \mathbb{E}_{Y|X=x}[Y | X = x] + y^2}_{J(y)}$$

$$\Leftrightarrow \forall x \in \mathcal{X}, f^*(x) = \arg \min_{y \in \mathbb{R}} J(y)$$

- Minimum obtenu lorsque la dérivée est nulle :

$$\frac{dJ(y)}{dy} = -2\mathbb{E}_{Y|X=x}[Y | X = x] + 2y = 0 \Rightarrow y = \mathbb{E}_{Y|X=x}[Y | X = x]$$

Décomposition de l'erreur

Par rapport au meilleur modèle de la classe f^*

- ▶ Soit f^* la meilleure fonction de \mathcal{F} : $f^* = \arg \min_{f \in \mathcal{F}} R(f)$
- ▶ $R(f_{\text{Bayes/reg}}) = \inf_f R(f)$: risque du modèle optimal
(sans obligation d'avoir $f_{\text{Bayes/reg}} \in \mathcal{F}$)
- ▶ Le risque d'un modèle \hat{f} sélectionné par un algorithme se décompose en

$$R(\hat{f}) - R(f_{\text{Bayes/reg}}) = \underbrace{[R(\hat{f}) - R(f^*)]}_{\text{erreur d'estimation}} + \underbrace{[R(f^*) - R(f_{\text{Bayes/reg}})]}_{\text{erreur d'approximation}}$$

- ▶ **Erreur d'approximation** : erreur due au choix de \mathcal{F}
 $(= 0 \text{ si } f_{\text{Bayes/reg}} \in \mathcal{F})$
- ▶ **Erreur d'estimation** : erreur dépendant des données

Décomposition de l'erreur

Par rapport au meilleur modèle de la classe f^*

- ▶ Soit f^* la meilleure fonction de \mathcal{F} : $f^* = \arg \min_{f \in \mathcal{F}} R(f)$
- ▶ $R(f_{\text{Bayes/reg}}) = \inf_f R(f)$: risque du modèle optimal
(sans obligation d'avoir $f_{\text{Bayes/reg}} \in \mathcal{F}$)
- ▶ Le risque d'un modèle \hat{f} sélectionné par un algorithme se décompose en

$$R(\hat{f}) - R(f_{\text{Bayes/reg}}) = \underbrace{[R(\hat{f}) - R(f^*)]}_{\text{erreur d'estimation}} + \underbrace{[R(f^*) - R(f_{\text{Bayes/reg}})]}_{\text{erreur d'approximation}}$$

- ▶ **Erreur d'approximation** : erreur due au choix de \mathcal{F}
 $(= 0 \text{ si } f_{\text{Bayes/reg}} \in \mathcal{F})$
- ▶ **Erreur d'estimation** : erreur dépendant des données

Dilemme

Plus la classe de fonctions \mathcal{F} est étendue, plus l'erreur d'approximation peut être faible, mais plus l'erreur d'estimation peut être grande
(il est plus difficile de trouver f^* dans un ensemble plus vaste de modèles)

Surapprentissage (*overfitting*)

ERM *a priori* efficace si le modèle optimal est dans \mathcal{F} , alors

Pourquoi restreindre la classe de fonctions \mathcal{F} ?

- ▶ Imaginons un problème de discrimination binaire déterministe avec
 $Y = f_{Bayes}(X)$

Surapprentissage (*overfitting*)

ERM *a priori* efficace si le modèle optimal est dans \mathcal{F} , alors

Pourquoi restreindre la classe de fonctions \mathcal{F} ?

- ▶ Imaginons un problème de discrimination binaire déterministe avec $Y = f_{Bayes}(X)$
- ▶ Sans contrainte sur \mathcal{F} , il est toujours possible de choisir f telle que pour une base d'apprentissage S :
 - ▶ $\forall (x_i, y_i) \in S, f(x_i) = y_i$
 - ▶ et $f(x) \neq f_{Bayes}(x)$ sur tous les autres $x \in \mathcal{X}$

(cas extrême d'apprentissage par cœur)

ce qui conduit à

$$R_{emp}(f) = 0, \quad \text{et} \quad R(f) = 1$$

Surapprentissage (*overfitting*)

ERM *a priori* efficace si le modèle optimal est dans \mathcal{F} , alors

Pourquoi restreindre la classe de fonctions \mathcal{F} ?

- ▶ Imaginons un problème de discrimination binaire déterministe avec $Y = f_{Bayes}(X)$
- ▶ Sans contrainte sur \mathcal{F} , il est toujours possible de choisir f telle que pour une base d'apprentissage S :
 - ▶ $\forall (x_i, y_i) \in S, f(x_i) = y_i$
 - ▶ et $f(x) \neq f_{Bayes}(x)$ sur tous les autres $x \in \mathcal{X}$

(cas extrême d'apprentissage par cœur)

ce qui conduit à

$$R_{emp}(f) = 0, \quad \text{et} \quad R(f) = 1$$

But : trouver le bon compromis conduisant à une erreur empirique faible sans surapprentissage

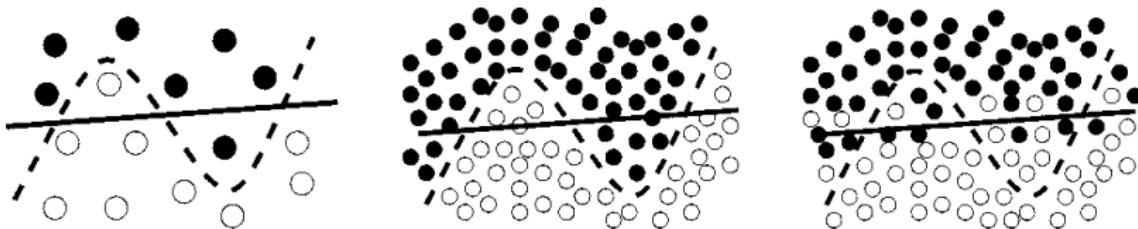
Surapprentissage (*overfitting*)

Le bon compromis

- ▶ Surapprentissage \approx apprentissage d'une fonction trop complexe
- ▶ Trouver le compromis entre l'erreur d'apprentissage faible et la simplicité du classifieur

Difficile à trouver

- ▶ Dépend du problème
- ▶ Dépend des données



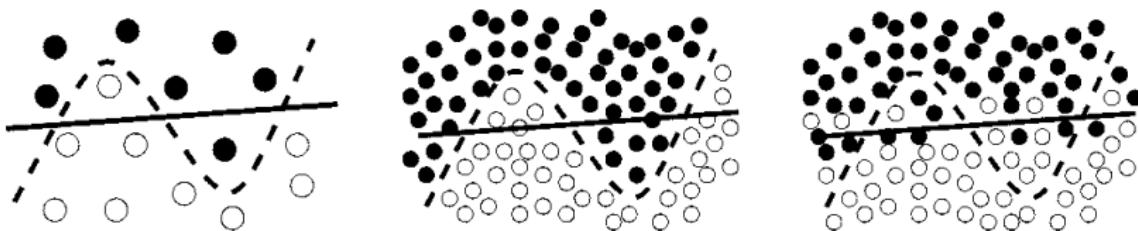
Surapprentissage (*overfitting*)

Le bon compromis

- ▶ Surapprentissage \approx apprentissage d'une fonction trop complexe
- ▶ Trouver le compromis entre l'erreur d'apprentissage faible et la simplicité du classifieur

Difficile à trouver

- ▶ Dépend du problème
- ▶ Dépend des données



En pratique : trouver le bon compromis = régler un (ou plusieurs) hyperparamètres

Régression linéaire

Restreindre la classe de fonctions

- ▶ Pour éviter le surapprentissage ; et pour faciliter l'apprentissage
- ▶ Classe des modèles linéaires :

$$\mathcal{F} = \left\{ f \mid f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \mathbf{w}^T \mathbf{x} = \sum_{j=1}^d w_j x_j, \mathbf{w} \in \mathbb{R}^d \right\}$$

Méthode des moindres carrés

- ▶ Application directe du principe ERM : minimiser l'EQM empirique

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

- ▶ Forme analytique de la solution :

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \text{avec } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

- ▶ Pour les modèles affines, \mathbf{x} contient une composante constante ($= 1$)

Méthode des moindres carrés

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 = \arg \min_{\mathbf{w} \in \mathbb{R}^d} J(\mathbf{w})$$

$$\text{avec } J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Gradient nul au minimum : $\nabla J(\hat{\mathbf{w}}) = \mathbf{0}$

$$J(\mathbf{w}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w}$$

$$\nabla J(\mathbf{w}) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\mathbf{w}$$

$$\nabla J(\hat{\mathbf{w}}) = \mathbf{0} \Rightarrow \mathbf{X}^T \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y} \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Sans notation matricielle :

$$J(\mathbf{w}) = \sum_{i=1}^m y_i^2 - 2y_i \mathbf{x}_i^T \mathbf{w} + \mathbf{w}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w}$$

$$\nabla J(\mathbf{w}) = -2 \sum_{i=1}^m y_i \mathbf{x}_i + 2 \left(\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w}$$

$$\nabla J(\hat{\mathbf{w}}) = \mathbf{0} \Rightarrow \left(\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \right) \hat{\mathbf{w}} = \sum_{i=1}^m y_i \mathbf{x}_i \Rightarrow \hat{\mathbf{w}} = \left(\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \sum_{i=1}^m y_i \mathbf{x}_i$$

Classification binaire et linéaire

Classificateurs binaires ($\mathcal{Y} = \{-1, +1\}$)

- ▶ Fonction de score à valeurs réelles $g : \mathcal{X} \rightarrow \mathbb{R}$
- ▶ Classificateur binaire : $f(\mathbf{x}) = \text{signe}(g(\mathbf{x}))$

Classificateurs linéaires

- ▶ A base de fonctions linéaires $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$:

$$\mathcal{F} = \left\{ f \mid f(\mathbf{x}) = \text{signe} \left(\mathbf{w}^T \mathbf{x} + b \right), \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

Apprentissage

- ▶ Dans le cas général : application directe du principe ERM "impossible" car la fonction de perte $\mathbb{I}(f(\mathbf{x}_i) \neq y_i)$ est non convexe et non différentiable
⇒ on considère souvent une approximation de la fonction de perte
- ▶ Dans le cas linéairement séparable ($\exists f \in \mathcal{F}, R_{emp}(f) = 0$) : problème de programmation linéaire

Séparabilité linéaire

Séparabilité linéaire pour $\mathcal{Y} = \{-1, +1\}$

- ▶ Un jeu de données S est linéairement séparable s'il existe un hyperplan séparant l'espace de représentation en deux parties ne contenant respectivement que des exemples de chacune des classes
- ▶ **Hyperplan** : surface $H = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^T \mathbf{x} + b = 0\}$
- ▶ $\mathbf{w} \in \mathbb{R}^d$: vecteur normal à l'hyperplan
- ▶ $b \in \mathbb{R}$: biais (décallage par rapport à l'origine)
suivant la dimension d , l'hyperplan peut être un point, une droite, un plan, ou un hyperplan
- ▶ Distance d'un point \mathbf{x} à l'hyperplan (pour $\|\mathbf{w}\| = 1$) :

$$d(\mathbf{x}, H) = |\mathbf{w}^T \mathbf{x} + b|$$

- ▶ Classification par un hyperplan de séparation :

$$f(\mathbf{x}) = \text{signe}(\mathbf{w}^T \mathbf{x} + b)$$

- ▶ S est linéairement séparable s'il existe $(\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$ tels que

$$\forall i \in \{1, \dots, m\}, \quad \text{signe}(\mathbf{w}^T \mathbf{x}_i + b) = y_i$$

Séparabilité linéaire et programmation linéaire

S est linéairement séparable s'il existe $(\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$ tels que

$$\forall i \in \{1, \dots, m\}, \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) > 0$$

Problème de programmation linéaire (LP)

Soit

$$\tilde{\mathbf{X}} = \begin{bmatrix} y_1 \mathbf{x}_1^T, & y_1 \\ \vdots \\ y_m \mathbf{x}_m^T, & y_m \end{bmatrix}$$

S est linéairement séparable si $\exists (\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$, tel que

$$\tilde{\mathbf{X}} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} > 0$$

ou encore, si $\exists (\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$, tel que

$$\tilde{\mathbf{X}} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \geq 1$$

Un des premiers algorithmes d'apprentissage

Le Perceptron (Rosenblatt, 1958)

- ▶ Le perceptron implémente la fonction

$$f(\mathbf{x}) = \begin{cases} +1, & \text{si } \theta^T \mathbf{x} \geq \theta_0 \\ -1, & \text{sinon} \end{cases}$$

$\theta \in \mathbb{R}^d$: vecteur des paramètres ("poids")

$\theta_0 \in \mathbb{R}$: paramètre de biais (seuil de décision)

Un des premiers algorithmes d'apprentissage

Le Perceptron (Rosenblatt, 1958)

- ▶ Le perceptron implémente la fonction

$$f(\mathbf{x}) = \begin{cases} +1, & \text{si } \boldsymbol{\theta}^T \mathbf{x} \geq \theta_0 \\ -1, & \text{sinon} \end{cases}$$

$\boldsymbol{\theta} \in \mathbb{R}^d$: vecteur des paramètres ("poids")

$\theta_0 \in \mathbb{R}$: paramètre de biais (seuil de décision)

Interprétation géométrique

- ▶ La classe de fonctions correspondante est la famille des classificateurs linéaires de \mathbb{R}^d :

$$\mathcal{F} = \left\{ f \mid f(\mathbf{x}) = \text{signe}(\boldsymbol{\theta}^T \mathbf{x} - \theta_0), \theta_0 \in \mathbb{R}, \boldsymbol{\theta} \in \mathbb{R}^d \right\}$$

- ▶ *Hyperplan* de séparation dans l'espace de représentation

Perceptron

Algorithme d'apprentissage du perceptron (cas sans biais : $\theta_0 = 0$)

Initialisation (aléatoire) des poids θ

Répéter jusqu'à convergence :

Pour i de 1 à m ,

$$\theta \leftarrow \begin{cases} \theta, & \text{si } f(\mathbf{x}_i) = y_i \\ \theta + y_i \mathbf{x}_i, & \text{sinon} \end{cases}$$

(chaque itération sur la base d'apprentissage est appelée une *epoch*)

Remarques

- ▶ Le terme de biais fixé à 0 implique que l'hyperplan séparateur passe par l'origine
- ▶ En pratique : choisir la séquence des exemples traités aléatoirement
- ▶ Interprétation : correction positive si f ne prend pas assez en compte l'exemple positif \mathbf{x}_i ou négative si l'influence de \mathbf{x}_i est trop grande

Perceptron

Algorithme d'apprentissage du perceptron (cas sans biais : $\theta_0 = 0$)

Initialisation (aléatoire) des poids θ

Répéter jusqu'à convergence :

Pour i de 1 à m ,

$$\theta \leftarrow \begin{cases} \theta, & \text{si } f(\mathbf{x}_i) = y_i \\ \theta + y_i \mathbf{x}_i, & \text{sinon} \end{cases}$$

(chaque itération sur la base d'apprentissage est appelée une *epoch*)

Remarques

- ▶ Le terme de biais fixé à 0 implique que l'hyperplan séparateur passe par l'origine
- ▶ En pratique : choisir la séquence des exemples traités aléatoirement
- ▶ Interprétation : correction positive si f ne prend pas assez en compte l'exemple positif \mathbf{x}_i ou négative si l'influence de \mathbf{x}_i est trop grande

Convergence

Si la base d'apprentissage S est linéairement séparable, l'algorithme du perceptron converge en un nombre fini d'itérations

Perceptron

Algorithme d'apprentissage du perceptron (avec biais)

$$\theta^T \mathbf{x} \geq \theta_0 \Leftrightarrow \theta^T \mathbf{x} - \theta_0 \geq 0 \Leftrightarrow [\theta^T \quad \theta_0] \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} \geq 0 \Leftrightarrow \tilde{\theta}^T \tilde{\mathbf{x}} \geq 0$$

Perceptron

Algorithme d'apprentissage du perceptron (avec biais)

$$\theta^T \mathbf{x} \geq \theta_0 \Leftrightarrow \theta^T \mathbf{x} - \theta_0 \geq 0 \Leftrightarrow [\theta^T \quad \theta_0] \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} \geq 0 \Leftrightarrow \tilde{\theta}^T \tilde{\mathbf{x}} \geq 0$$

Algorithme similaire au cas sans biais en dimension $d + 1$:

Initialisation (aléatoire) des poids $\tilde{\theta} = \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix}$

Répéter jusqu'à convergence :

Pour i de 1 à m ,

$$\tilde{\theta} \leftarrow \begin{cases} \tilde{\theta}, & \text{si } f(\mathbf{x}_i) = y_i \\ \tilde{\theta} + y_i \tilde{\mathbf{x}}_i, & \text{sinon} \end{cases}$$

$$\Leftrightarrow \theta \leftarrow \begin{cases} \theta, & \text{si } f(\mathbf{x}_i) = y_i \\ \theta + y_i \mathbf{x}_i, & \text{sinon} \end{cases} \quad \text{et} \quad \theta_0 \leftarrow \begin{cases} \theta_0, & \text{si } f(\mathbf{x}_i) = y_i \\ \theta_0 - y_i, & \text{sinon} \end{cases}$$

Perceptron

Algorithme d'apprentissage du perceptron (avec biais)

$$\theta^T \mathbf{x} \geq \theta_0 \Leftrightarrow \theta^T \mathbf{x} - \theta_0 \geq 0 \Leftrightarrow [\theta^T \quad \theta_0] \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} \geq 0 \Leftrightarrow \tilde{\theta}^T \tilde{\mathbf{x}} \geq 0$$

Algorithme similaire au cas sans biais en dimension $d + 1$:

Initialisation (aléatoire) des poids $\tilde{\theta} = \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix}$

Répéter jusqu'à convergence :

Pour i de 1 à m ,

$$\tilde{\theta} \leftarrow \begin{cases} \tilde{\theta}, & \text{si } f(\mathbf{x}_i) = y_i \\ \tilde{\theta} + y_i \tilde{\mathbf{x}}_i, & \text{sinon} \end{cases}$$

$$\Leftrightarrow \theta \leftarrow \begin{cases} \theta, & \text{si } f(\mathbf{x}_i) = y_i \\ \theta + y_i \mathbf{x}_i, & \text{sinon} \end{cases} \quad \text{et} \quad \theta_0 \leftarrow \begin{cases} \theta_0, & \text{si } f(\mathbf{x}_i) = y_i \\ \theta_0 - y_i, & \text{sinon} \end{cases}$$

En pratique...

Avec un taux d'apprentissage $\mu \in (0, 1)$: $\tilde{\theta} \leftarrow \begin{cases} \tilde{\theta}, & \text{si } f(\mathbf{x}_i) = y_i \\ \tilde{\theta} + \mu y_i \tilde{\mathbf{x}}_i, & \text{sinon} \end{cases}$

Algorithme des k plus proches voisins

Intuition

- ▶ Classer l'exemple x selon la classe de l'exemple de la base d'apprentissage le plus proche
⇒ *très sensible aux données*
- ▶ "Vote majoritaire" sur un ensemble d'exemples proches de x
⇒ *sensibilité aux données atténuée*

Algorithme des k plus proches voisins

Intuition

- ▶ Classer l'exemple x selon la classe de l'exemple de la base d'apprentissage le plus proche
⇒ *très sensible aux données*
- ▶ "Vote majoritaire" sur un ensemble d'exemples proches de x
⇒ *sensibilité aux données atténuée*

Algorithme pour $\mathcal{Y} = \{1, \dots, Q\}$

$$f(x) = \arg \max_{j=1, \dots, Q} \sum_{x_i \in \mathcal{V}_k(x)} \mathbb{I}(y_i = j)$$

- ▶ $\mathcal{V}_k(x)$: k -voisinage de x (k plus proches exemples de x)

Algorithme des k plus proches voisins

Intuition

- ▶ Classer l'exemple x selon la classe de l'exemple de la base d'apprentissage le plus proche
⇒ *très sensible aux données*
- ▶ "Vote majoritaire" sur un ensemble d'exemples proches de x
⇒ *sensibilité aux données atténuée*

Algorithme pour $\mathcal{Y} = \{1, \dots, Q\}$

$$f(x) = \arg \max_{j=1, \dots, Q} \sum_{x_i \in \mathcal{V}_k(x)} \mathbb{I}(y_i = j)$$

- ▶ $\mathcal{V}_k(x)$: k -voisinage de x (k plus proches exemples de x)

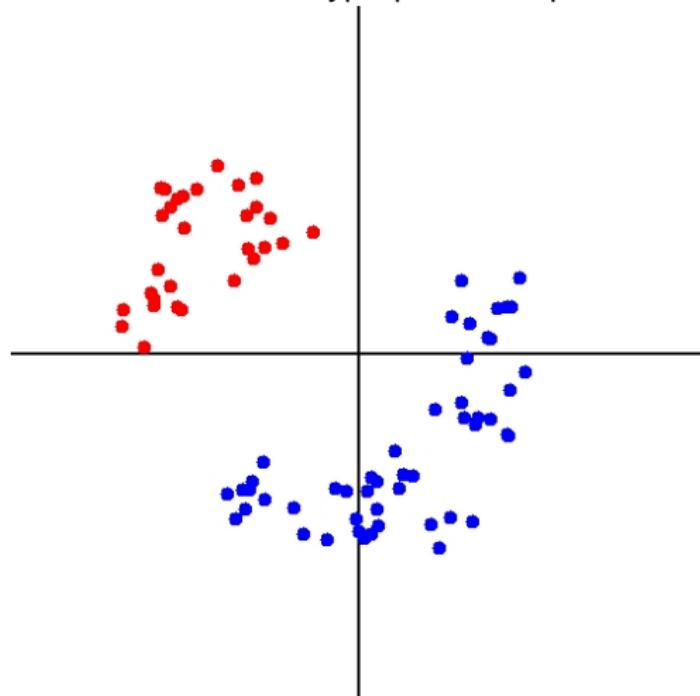
Remarques

- + Pas de phase d'apprentissage autre que "stocker tous les exemples"
- + Cas multi-classe équivalent au cas bi-classe
- + Pas d'hypothèse sur la séparabilité linéaire des données
- Nécessité de retenir toute la base d'apprentissage en mémoire
- Nécessité de calculer toutes les distances à chaque prédiction

Hyperplan de séparation optimal

Quel est le meilleur classifieur linéaire $f(\mathbf{x}) = \text{signe}(\mathbf{w}^T \mathbf{x} + b)$?

Quel est le meilleur hyperplan de séparation $H = \{\mathbf{x} | \mathbf{w}^T \mathbf{x} + b = 0\}$?



- ▶ La **marge** d'un classifieur linéaire f par rapport à un jeu de données S est la plus petite distance d'un exemple \mathbf{x}_i à l'hyperplan séparateur H correspondant

Hyperplan de séparation optimal

Marge d'un classifieur

- Marge Δ d'un classifieur linéaire par rapport à un jeu de données S :

$$\Delta = \min_{\mathbf{x}_i \in S} dist(\mathbf{x}_i, H) = \min_{\mathbf{x}_i \in S} \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

Hyperplan de séparation optimal

Marge d'un classifieur

- Marge Δ d'un classifieur linéaire par rapport à un jeu de données S :

$$\Delta = \min_{\mathbf{x}_i \in S} dist(\mathbf{x}_i, H) = \min_{\mathbf{x}_i \in S} \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

Apprendre en cherchant à maximiser la marge

- Problème :
 - H est insensible à un changement d'échelle des paramètres (\mathbf{w}, b) :
 $\mathbf{w}^T \mathbf{x} + b = 0 \Leftrightarrow 2\mathbf{w}^T \mathbf{x} + 2b = 0$
 - Il existe une infinité de paramètres (\mathbf{w}, b) solution
- Fixer $\|\mathbf{w}\| = 1$? \rightarrow pas évident d'un point de vue pratique
- **Hyperplan canonique** : pour \mathbf{x}_k sur la marge
 $(\mathbf{x}_k = \arg \min_{\mathbf{x}_i \in S} dist(\mathbf{x}_i, H))$,

$$fixer \quad |\mathbf{w}^T \mathbf{x}_k + b| = 1 \quad \Rightarrow \quad \Delta = \frac{1}{\|\mathbf{w}\|}$$

- Maximiser la marge revient à minimiser $\|\mathbf{w}\|$

Machine à Vecteurs Support (SVM) à marge dure

Hard-margin Support Vector Machine

- ▶ **But** : maximiser la marge $\Rightarrow \min \|w\|$
- ▶ **Contrainte de marge dure** : tous les exemples doivent être bien classés

$$\forall i, \begin{cases} w^T x_i + b \geq 0, & \text{si } y_i = +1 \\ w^T x_i + b < 0, & \text{si } y_i = -1 \end{cases}$$

et en dehors de la marge :

$$\forall i, \text{dist}(x_i, H) = \frac{|w^T x_i + b|}{\|w\|} \geq \Delta = \frac{1}{\|w\|} \quad \Rightarrow \quad |w^T x_i + b| \geq 1$$

Machine à Vecteurs Support (SVM) à marge dure

Hard-margin Support Vector Machine

- ▶ **But** : maximiser la marge $\Rightarrow \min \|\mathbf{w}\|$
- ▶ **Contrainte de marge dure** : tous les exemples doivent être bien classés

$$\forall i, \begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 0, & \text{si } y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b < 0, & \text{si } y_i = -1 \end{cases}$$

et en dehors de la marge :

$$\forall i, \text{dist}(\mathbf{x}_i, H) = \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|} \geq \Delta = \frac{1}{\|\mathbf{w}\|} \quad \Rightarrow \quad |\mathbf{w}^T \mathbf{x}_i + b| \geq 1$$

- ▶ Formulation de l'apprentissage :

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.c. } & \begin{cases} y_1(\mathbf{w}^T \mathbf{x}_1 + b) \geq 1 \\ \vdots \\ y_m(\mathbf{w}^T \mathbf{x}_m + b) \geq 1 \end{cases} \end{aligned}$$

- ▶ Problème de programmation quadratique convexe : solution unique, algorithmes efficaces disponibles

Machine à Vecteurs Support (SVM) à marge dure

Hard-margin Support Vector Machine

- ▶ **But** : maximiser la marge $\Rightarrow \min \|\mathbf{w}\|$
- ▶ **Contrainte de marge dure** : tous les exemples doivent être bien classés

$$\forall i, \begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 0, & \text{si } y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b < 0, & \text{si } y_i = -1 \end{cases}$$

et en dehors de la marge :

$$\forall i, \text{dist}(\mathbf{x}_i, H) = \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|} \geq \Delta = \frac{1}{\|\mathbf{w}\|} \quad \Rightarrow \quad |\mathbf{w}^T \mathbf{x}_i + b| \geq 1$$

- ▶ Formulation de l'apprentissage :

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.c. } & \begin{cases} y_1(\mathbf{w}^T \mathbf{x}_1 + b) \geq 1 \\ \vdots \\ y_m(\mathbf{w}^T \mathbf{x}_m + b) \geq 1 \end{cases} \end{aligned}$$

- ▶ Problème de programmation quadratique convexe : solution unique, algorithmes efficaces disponibles
- ▶ uniquement pour des problèmes linéairement séparables

Machine à Vecteurs Support (SVM) à marge dure

Propriétés

- ▶ Les exemples \mathbf{x}_i sur la marge (avec $|\mathbf{w}^T \mathbf{x}_i + b| = 1$) sont appelés *vecteurs support*
- ▶ Ils "supportent" la solution et contiennent toute l'information nécessaire à la résolution du problème
- ▶ Déplacer ou supprimer un exemple x_i en dehors de la marge ne change pas la solution

- ▶ La formulation duale du problème d'optimisation permet de s'affranchir du problème de la dimension des données

Machine à Vecteurs Support (SVM) à marge dure

Propriétés

- ▶ Les exemples \mathbf{x}_i sur la marge (avec $|\mathbf{w}^T \mathbf{x}_i + b| = 1$) sont appelés *vecteurs support*
- ▶ Ils "supportent" la solution et contiennent toute l'information nécessaire à la résolution du problème
- ▶ Déplacer ou supprimer un exemple x_i en dehors de la marge ne change pas la solution

- ▶ La formulation duale du problème d'optimisation permet de s'affranchir du problème de la dimension des données

Cas non séparable

- ▶ SVM à marge douce...

Machine à Vecteurs Support (SVM)

SVM à marge douce

- ▶ Relâcher les contraintes pour rendre le problème réalisable
- ▶ Optimiser un compromis entre la quantité d'erreur et la marge

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi_i \in \mathbb{R}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.c. } & \left\{ \begin{array}{l} y_1(\mathbf{w}^T \mathbf{x}_1 + b) \geq 1 - \xi_1 \\ \vdots \\ y_m(\mathbf{w}^T \mathbf{x}_m + b) \geq 1 - \xi_m \\ \xi_1 \geq 0, \dots, \xi_m \geq 0 \end{array} \right. \end{aligned}$$

$C > 0$: hyperparamètre permettant de régler le compromis

- ▶ C petit favorise la maximisation de la marge et conduit à plus d'erreurs
- ▶ C grand impose peu d'erreurs et conduit à une marge plus faible
- ▶ Toujours un problème de programmation quadratique convexe
- ▶ Utilisable sur des données non linéairement séparables
- ▶ A favoriser dans tous les cas

Plan

Qu'est-ce que l'IA ?

Résolution de problèmes

Stratégies de jeux

Modélisation de l'incertain

Apprentissage supervisé

Classification multi-classe

Méthodes non linéaires

Méthodes à noyaux

Problèmes multi-classes : $\mathcal{Y} = \{1, \dots, Q\}$

- ▶ Problème de classification avec un nombre de classes $Q \geq 3$
- ▶ Fonction de perte de classification standard : $\ell(f, X, Y) = \mathbb{I}(f(X) \neq Y)$

Classifieur multi-classe optimal (classifieur de Bayes)

- ▶ Prédire la classe la plus probable pour l'exemple :

$$f_{Bayes}(x) = \arg \max_{y \in \{1, \dots, Q\}} P(Y = y | X = x)$$

Problèmes multi-classes : $\mathcal{Y} = \{1, \dots, Q\}$

- ▶ Problème de classification avec un nombre de classes $Q \geq 3$
- ▶ Fonction de perte de classification standard : $\ell(f, X, Y) = \mathbb{I}(f(X) \neq Y)$

Classifieur multi-classe optimal (classifieur de Bayes)

- ▶ Prédire la classe la plus probable pour l'exemple :

$$f_{Bayes}(x) = \arg \max_{y \in \{1, \dots, Q\}} P(Y = y | X = x)$$

Risque de Bayes

- ▶ Cas déterministe : $P(Y = y | X = x) \in \{0, 1\}$
$$\Rightarrow P(Y = f_{Bayes}(X) | X = x) = 1 \quad \Rightarrow \quad R(f_{Bayes}) = 0$$

Problèmes multi-classes : $\mathcal{Y} = \{1, \dots, Q\}$

- ▶ Problème de classification avec un nombre de classes $Q \geq 3$
- ▶ Fonction de perte de classification standard : $\ell(f, X, Y) = \mathbb{I}(f(X) \neq Y)$

Classifieur multi-classe optimal (classifieur de Bayes)

- ▶ Prédire la classe la plus probable pour l'exemple :

$$f_{Bayes}(x) = \arg \max_{y \in \{1, \dots, Q\}} P(Y = y | X = x)$$

Risque de Bayes

- ▶ Cas déterministe : $P(Y = y | X = x) \in \{0, 1\}$

$$\Rightarrow P(Y = f_{Bayes}(X) | X = x) = 1 \quad \Rightarrow \quad R(f_{Bayes}) = 0$$

- ▶ Cas stochastique : $P(Y = y | X = x) \in [0, 1]$

$$P(Y \neq f_{Bayes}(X) | X = x) = \sum_{y \neq f_{Bayes}(x)} P(Y = y | X = x) = 1 - P(Y = f_{Bayes}(X) | X = x)$$

$$\Rightarrow R(f_{Bayes}) = \mathbb{E}_X P(Y \neq f_{Bayes}(X) | X = x) = \mathbb{E}_X [1 - \max_{y \in \mathcal{Y}} P(Y = y | X = x)]$$

Problèmes multi-classes

Méthodes de décomposition

Construction du classifieur multi-classe à partir d'un ensemble de classificateurs binaires

- ▶ **Méthode un-contre-tous** : Q classificateurs binaires

chaque classificateur binaire $f_j(x) = \text{signe}(g_j(x))$ est chargé de séparer une classe de l'ensemble des autres

$$f(x) = \arg \max_{j \in \{1, \dots, Q\}} g_j(x)$$

- ▶ **Méthode un-contre-un** : $Q(Q - 1)/2$ classificateurs binaires

un classificateur binaire $f_{j/k}(x) \in \{j, k\}$ pour séparer chaque couple de classes (j, k)

$f(x) =$ classe avec le "vote majoritaire"

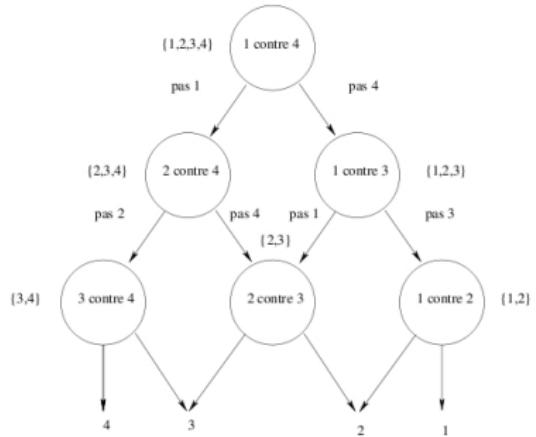
$$= \arg \max_{y \in \{1, \dots, Q\}} \sum_{1 \leq j < k \leq Q} \mathbb{I}(f_{j/k}(x) = y)$$

Problèmes multi-classes

Méthodes de décomposition (suite)

► Méthodes basées sur un arbre de décision :

pour limiter le nombre de calculs en prédiction avec une approche un-contre-un

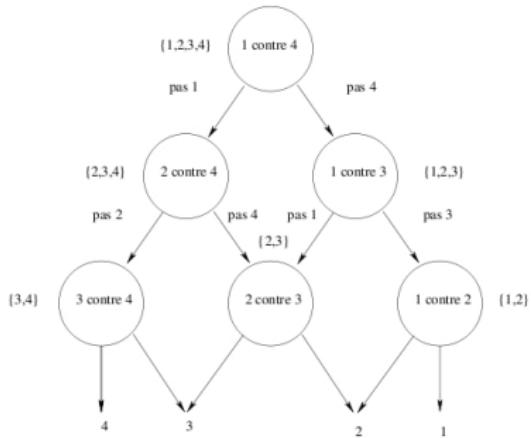


Problèmes multi-classes

Méthodes de décomposition (suite)

► Méthodes basées sur un arbre de décision :

pour limiter le nombre de calculs en prédiction avec une approche un-contre-un



Méthodes directes

- Classifieurs à "sorties" multiples (réseaux de neurones, M-SVM, etc.)
- k plus proches voisins, classifieur naïf de Bayes

Plan

Introduction

Apprentissage automatique

Liens avec la théorie de l'estimation

Rappels de proba/stats

Apprentissage supervisé

Modélisation probabiliste de l'apprentissage

Surapprentissage

Régression linéaire

Classification linéaire

Algorithmes classiques en classification

Classification multi-classe

Régularisation

Méthodes non linéaires

Méthodes à noyaux

Théorie statistique de l'apprentissage

Estimation du risque

Bornes de généralisation

Apprentissage non supervisé

Régularisation

Contrôler la complexité du modèle pendant l'apprentissage

- ▶ En minimisant un compromis entre
 - ▶ un terme d'erreur (ou d'attaché aux données)
 - ▶ et un terme de régularisation, $\mathcal{R}(f)$, pénalisant les fonctions complexes

$$\min_{f \in \mathcal{F}} \sum_{i=1}^m \ell(f, x_i, y_i) + \lambda \mathcal{R}(f)$$

λ : hyperparamètre pondérant la régularisation
(et permettant de régler le compromis)

Exemple des SVM

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \underbrace{\sum_{i=1}^m \max\{0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)\}}_{\text{terme d'erreur}} + \lambda \underbrace{\|\mathbf{w}\|^2}_{\text{terme de régularisation}}$$

ℓ : fonction de perte "charnière" (*hinge loss*)

λ joue un rôle similaire à $1/2C$ dans la formulation SVM précédente

Régression régularisée

- ▶ Utile en grande dimension
- ▶ Indispensable dès que $d > m$:
 $d > m \Rightarrow \text{rank}(\mathbf{X}^T \mathbf{X}) < d \Rightarrow$ pas de solution unique des moindres carrés

Régression régularisée

- ▶ Utile en grande dimension
- ▶ Indispensable dès que $d > m$:
 $d > m \Rightarrow \text{rank}(\mathbf{X}^T \mathbf{X}) < d \Rightarrow$ pas de solution unique des moindres carrés

Régression "ridge" / moindres carrés régularisés

- ▶ Régression linéaire, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, avec
 - ▶ perte quadratique : $\ell(f, \mathbf{x}_i, y_i) = (y_i - f(\mathbf{x}_i))^2$
 - ▶ régularisation de type ℓ_2 : $\mathcal{R}(f) = \|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w}$
- ▶ Solution analytique :

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Régression régularisée

- ▶ Utile en grande dimension
- ▶ Indispensable dès que $d > m$:
 $d > m \Rightarrow \text{rank}(\mathbf{X}^T \mathbf{X}) < d \Rightarrow$ pas de solution unique des moindres carrés

Régression "ridge" / moindres carrés régularisés

- ▶ Régression linéaire, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, avec
 - ▶ perte quadratique : $\ell(f, \mathbf{x}_i, y_i) = (y_i - f(\mathbf{x}_i))^2$
 - ▶ régularisation de type ℓ_2 : $\mathcal{R}(f) = \|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w}$

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_2^2 = \|\mathbf{y} - \mathbf{Xw}\|_2^2 + \lambda \mathbf{w}^T \mathbf{w}$$

- ▶ Solution analytique :

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Régression LASSO

- ▶ Idem avec régularisation de type ℓ_1 : $\mathcal{R}(f) = \|\mathbf{w}\|_1 = \sum_{k=1}^d |w_k|$
- ▶ Favorise la parcimonie ($w_k = 0$ pour un certain nombre de paramètres)
- ▶ Solution analytique si \mathbf{X} est orthonormée
- ▶ En général : problème de programmation quadratique (solveurs efficaces disponibles)

Plan

Introduction

Apprentissage automatique

Liens avec la théorie de l'estimation

Rappels de proba/stats

Apprentissage supervisé

Modélisation probabiliste de l'apprentissage

Surapprentissage

Régression linéaire

Classification linéaire

Algorithmes classiques en classification

Classification multi-classe

Régularisation

Méthodes non linéaires

Méthodes à noyaux

Théorie statistique de l'apprentissage

Estimation du risque

Bornes de généralisation

Apprentissage non supervisé

Régression et classification non linéaires

Approche directe par projection non linéaire

- ▶ Projeter les données dans un nouvel espace de représentation où le problème devient linéaire
- ▶ Exemple :

$$f(x) = w_1 \sin(x) + w_2 x^2 - w_3 x^5 = \mathbf{w}^T \begin{bmatrix} \sin(x) \\ x^2 \\ -x^5 \end{bmatrix} = \mathbf{w}^T \phi(x)$$

- ▶ Apprendre le modèle non linéaire f revient à estimer les paramètres \mathbf{w} du modèle linéaire en $\phi(x)$
 1. Choisir une projection non linéaire $\phi : \mathcal{X} \rightarrow \mathcal{X}_\phi$
 2. Projeter toutes les données : $\Phi_i = \phi(\mathbf{x}_i)$, $i = 1, \dots, m$
 3. Appliquer une méthode linéaire aux exemples (Φ_i, y_i) au lieu de (\mathbf{x}_i, y_i)

Régression et classification non linéaires

Approche directe par projection non linéaire

- ▶ Projeter les données dans un nouvel espace de représentation où le problème devient linéaire
- ▶ Exemple :

$$f(x) = w_1 \sin(x) + w_2 x^2 - w_3 x^5 = \mathbf{w}^T \begin{bmatrix} \sin(x) \\ x^2 \\ -x^5 \end{bmatrix} = \mathbf{w}^T \phi(x)$$

- ▶ Apprendre le modèle non linéaire f revient à estimer les paramètres \mathbf{w} du modèle linéaire en $\phi(x)$
 1. Choisir une projection non linéaire $\phi : \mathcal{X} \rightarrow \mathcal{X}_\phi$
 2. Projeter toutes les données : $\Phi_i = \phi(\mathbf{x}_i)$, $i = 1, \dots, m$
 3. Appliquer une méthode linéaire aux exemples (Φ_i, y_i) au lieu de (\mathbf{x}_i, y_i)
- ▶ En général, la dimension de l'espace de représentation \mathcal{X}_ϕ augmente par rapport à celle de \mathcal{X}

Régression et classification non linéaires

Approche directe par projection non linéaire

- ▶ Projeter les données dans un nouvel espace de représentation où le problème devient linéaire
- ▶ Exemple :

$$f(x) = w_1 \sin(x) + w_2 x^2 - w_3 x^5 = \mathbf{w}^T \begin{bmatrix} \sin(x) \\ x^2 \\ -x^5 \end{bmatrix} = \mathbf{w}^T \phi(x)$$

- ▶ Apprendre le modèle non linéaire f revient à estimer les paramètres \mathbf{w} du modèle linéaire en $\phi(x)$
 1. Choisir une projection non linéaire $\phi : \mathcal{X} \rightarrow \mathcal{X}_\phi$
 2. Projeter toutes les données : $\Phi_i = \phi(\mathbf{x}_i)$, $i = 1, \dots, m$
 3. Appliquer une méthode linéaire aux exemples (Φ_i, y_i) au lieu de (\mathbf{x}_i, y_i)
- ▶ En général, la dimension de l'espace de représentation \mathcal{X}_ϕ augmente par rapport à celle de \mathcal{X}
- ▶ En augmentant suffisamment la dimension, il est possible d'approcher n'importe quelle non-linéarité

Régression et classification non linéaires

Approche directe par projection non linéaire

- ▶ Projeter les données dans un nouvel espace de représentation où le problème devient linéaire
- ▶ Exemple :

$$f(x) = w_1 \sin(x) + w_2 x^2 - w_3 x^5 = \mathbf{w}^T \begin{bmatrix} \sin(x) \\ x^2 \\ -x^5 \end{bmatrix} = \mathbf{w}^T \phi(x)$$

- ▶ Apprendre le modèle non linéaire f revient à estimer les paramètres \mathbf{w} du modèle linéaire en $\phi(x)$
 1. Choisir une projection non linéaire $\phi : \mathcal{X} \rightarrow \mathcal{X}_\phi$
 2. Projeter toutes les données : $\Phi_i = \phi(\mathbf{x}_i)$, $i = 1, \dots, m$
 3. Appliquer une méthode linéaire aux exemples (Φ_i, y_i) au lieu de (\mathbf{x}_i, y_i)
- ▶ En général, la dimension de l'espace de représentation \mathcal{X}_ϕ augmente par rapport à celle de \mathcal{X}
- ▶ En augmentant suffisamment la dimension, il est possible d'approcher n'importe quelle non-linéarité
- ▶ La régularisation (ou le contrôle de la complexité) devient cruciale

Plan

Introduction

Apprentissage automatique

Liens avec la théorie de l'estimation

Rappels de proba/stats

Apprentissage supervisé

Modélisation probabiliste de l'apprentissage

Surapprentissage

Régression linéaire

Classification linéaire

Algorithmes classiques en classification

Classification multi-classe

Régularisation

Méthodes non linéaires

Méthodes à noyaux

Théorie statistique de l'apprentissage

Estimation du risque

Bornes de généralisation

Apprentissage non supervisé

Méthodes à noyaux

Idée générale : Calculer des produits scalaires en grande dimension

- ▶ Permet de traiter des cas non linéaires complexes
 - ▶ Demande beaucoup de temps de calcul (et d'exemples)
- ⇒ on souhaiterait calculer $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ plus facilement

Exemple du noyau polynomial

Soit $\mathcal{X} = \mathbb{R}^2$ et

$$\phi : \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mapsto \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix},$$

Alors, le produit scalaire dans $\mathcal{X}_\phi \subset \mathbb{R}^3$ est

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} x'_1^2 \\ \sqrt{2}x'_1x'_2 \\ x'_2^2 \end{bmatrix} \right\rangle = x_1^2 x'^2_1 + 2x_1 x_2 x'_1 x'_2 + x_2^2 x'^2_2 \\ &= (x_1 x'_1 + x_2 x'_2)^2 \\ &= \langle \mathbf{x}, \mathbf{x}' \rangle^2\end{aligned}$$

On peut définir une fonction $K(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle^2$ qui calcule les produits scalaires dans \mathbb{R}^3 avec un produit scalaire dans \mathbb{R}^2 et une multiplication

Méthodes à noyaux

Noyaux définis positifs (*kernel functions*)

- Un noyau K est une fonction symétrique $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ définie positive, c'après que $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$ et

$$\forall n \in \mathbb{N}, \forall \{\mathbf{x}_i\}_{i=1}^n \in \mathcal{X}^n, \forall \{\alpha_i\}_{i=1}^n \in \mathbb{R}^n, \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

- Tout noyau K correspond à un produit scalaire dans un certain espace : il existe ϕ telle que

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

Astuce du noyau (*kernel trick*)

- Si le modèle et l'algorithme d'apprentissage peuvent se formuler uniquement à partir de produits scalaires entre vecteurs \mathbf{x}_i
- Alors il suffit de remplacer tous les $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ par $K(\mathbf{x}_i, \mathbf{x}_j)$ pour obtenir un modèle non linéaire
- Sans avoir besoin de calculer des $\phi(\mathbf{x})$ en grande dimension

Méthodes à noyaux

Noyaux usuels

- ▶ Noyau linéaire $K(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$
 ϕ = identité
- ▶ Noyau polynomial $K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle)^D$ ou $(1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^D$
 $\dim(\mathcal{X}_\phi) = \binom{d+D}{d} = \binom{D+d}{D} \geq (1 + d/D)^D$
- ▶ Noyau gaussien $K(\mathbf{x}, \mathbf{x}') = \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}\right)$
 $\dim(\mathcal{X}_\phi) = \infty$

Construction de nouveaux noyaux (par ex. dédiés à une application)

Si K_1 et K_2 sont des noyaux valides

- ▶ $K(\mathbf{x}, \mathbf{x}') = aK_1(\mathbf{x}, \mathbf{x}')$ (avec $a \geq 0$)
- ▶ $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + a$ (avec $a \geq 0$)
- ▶ $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$
- ▶ $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$
- ▶ $K(\mathbf{x}, \mathbf{x}') = \exp(K_1(\mathbf{x}, \mathbf{x}'))$

sont aussi des noyaux valides

SVM non linéaire

Formulation primale des SVM

- ▶ Modèle : $f(\mathbf{x}) = \text{signe}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$
- ▶ Apprentissage : programmation quadratique en dimension $d + m + 1$

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi_i \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$\begin{aligned} \text{s.c. } & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

SVM non linéaire

Formulation primale des SVM

- ▶ Modèle : $f(\mathbf{x}) = \text{signe}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$
- ▶ Apprentissage : programmation quadratique en dimension $d + m + 1$

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi_i \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$\begin{aligned} \text{s.c. } & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

Formulation duale des SVM

- ▶ Modèle : $f(\mathbf{x}) = \text{signe} \left(\sum_{i=1}^m \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b \right)$
- ▶ Apprentissage : programmation quadratique en dimension m

$$\max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^m \alpha_i$$

$$\text{s.c. } \sum_{i=1}^m \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C$$

- ▶ Vecteurs support : \mathbf{x}_i pour $\alpha_i \neq 0$

Formulation duale des SVM : démonstration

Ecriture du lagrangien

Formulation primale :

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi_i \in \mathbb{R}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$s.c. \quad y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

Formulation duale des SVM : démonstration

Ecriture du lagrangien

Formulation primale :

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi_i \in \mathbb{R}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{s.c. } y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i \Leftrightarrow y_i(\langle w, x_i \rangle + b) - 1 + \xi_i \geq 0 \\ \xi_i \geq 0$$

Lagrangien :

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(\langle w, x_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i$$

$\alpha_i \geq 0$ et $\mu_i \geq 0$: multiplicateurs de Lagrange (variables duales)

Formulation duale des SVM : démonstration

Obtention du lagrangien dual par minimisation du lagrangien

Lagrangien ($\alpha_i \geq 0$ et $\mu_i \geq 0$: multiplicateurs de Lagrange) :

$$L(\mathbf{w}, b, \xi, \alpha, \mu) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i$$

Formulation duale des SVM : démonstration

Obtention du lagrangien dual par minimisation du lagrangien

Lagrangien ($\alpha_i \geq 0$ et $\mu_i \geq 0$: multiplicateurs de Lagrange) :

$$L(\mathbf{w}, b, \xi, \alpha, \mu) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i$$

Point selle du Lagrangien : dérivées nulles par rapport aux variables primales :

$$\frac{\partial L(\mathbf{w}, b, \xi, \alpha, \mu)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = \mathbf{0} \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L(\mathbf{w}, b, \xi, \alpha, \mu)}{\partial b} = - \sum_{i=1}^m \alpha_i y_i = 0 \quad \Rightarrow \quad \sum_{i=1}^m \alpha_i y_i = 0$$

$$\frac{\partial L(\mathbf{w}, b, \xi, \alpha, \mu)}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad \Rightarrow \quad \mu_i = C - \alpha_i \quad \Rightarrow \quad 0 \leq \alpha_i \leq C$$

Formulation duale des SVM : démonstration

Obtention du lagrangien dual par minimisation du lagrangien

Lagrangien ($\alpha_i \geq 0$ et $\mu_i \geq 0$: multiplicateurs de Lagrange) :

$$L(\mathbf{w}, b, \xi, \alpha, \mu) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i$$

Point selle du lagrangien : dérivées nulles par rapport aux variables primales :

$$\frac{\partial L(\mathbf{w}, b, \xi, \alpha, \mu)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = \mathbf{0} \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L(\mathbf{w}, b, \xi, \alpha, \mu)}{\partial b} = - \sum_{i=1}^m \alpha_i y_i = 0 \quad \Rightarrow \quad \sum_{i=1}^m \alpha_i y_i = 0$$

$$\frac{\partial L(\mathbf{w}, b, \xi, \alpha, \mu)}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad \Rightarrow \quad \mu_i = C - \alpha_i \quad \Rightarrow \quad 0 \leq \alpha_i \leq C$$

Lagrangien dual : $L_D(\alpha, \mu) = \inf_{\mathbf{w}, b, \xi} L(\mathbf{w}, b, \xi, \alpha, \mu)$

$$\begin{aligned} L_D(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 + \xi_i \right] \\ &\quad - \sum_{i=1}^m (C - \alpha_i) \xi_i \end{aligned}$$

Formulation duale des SVM : démonstration

$$\begin{aligned}L_D(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 + \xi_i \right] \\&\quad - \sum_{i=1}^m (C - \alpha_i) \xi_i \\&= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 \right]\end{aligned}$$

Formulation duale des SVM : démonstration

$$\begin{aligned} L_D(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 + \xi_i \right] \\ &\quad - \sum_{i=1}^m (C - \alpha_i) \xi_i \\ &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 \right] \end{aligned}$$

On a

$$\sum_{i=1}^m \alpha_i y_i \left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 = \left\langle \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j \right\rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Formulation duale des SVM : démonstration

$$\begin{aligned} L_D(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 + \xi_i \right] \\ &\quad - \sum_{i=1}^m (C - \alpha_i) \xi_i \\ &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 \right] \end{aligned}$$

On a

$$\sum_{i=1}^m \alpha_i y_i \left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 = \left\langle \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j \right\rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Donc

$$L_D(\alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \dots$$

Formulation duale des SVM : démonstration

$$\begin{aligned}
 L_D(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 + \xi_i \right] \\
 &\quad - \sum_{i=1}^m (C - \alpha_i) \xi_i \\
 &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 \right]
 \end{aligned}$$

On a

$$\begin{aligned}
 \sum_{i=1}^m \alpha_i y_i \left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle &= \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\
 \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 &= \left\langle \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j \right\rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle
 \end{aligned}$$

Donc

$$L_D(\alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i$$

Formulation duale des SVM : démonstration

$$\begin{aligned}
 L_D(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 + \xi_i \right] \\
 &\quad - \sum_{i=1}^m (C - \alpha_i) \xi_i \\
 &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \alpha_i \left[y_i \left(\left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + b \right) - 1 \right]
 \end{aligned}$$

On a

$$\begin{aligned}
 \sum_{i=1}^m \alpha_i y_i \left\langle \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle &= \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\
 \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 &= \left\langle \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j \right\rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle
 \end{aligned}$$

Donc

$$\begin{aligned}
 L_D(\alpha) &= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - b \sum_{i=1}^m \cancel{\alpha_i y_i} + \sum_{i=1}^m \alpha_i \\
 &= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^m \alpha_i \quad \left(\text{car } \sum_{i=1}^m \alpha_i y_i = 0 \right)
 \end{aligned}$$

SVM non linéaire

Formulation duale des SVM

- Modèle : $f(\mathbf{x}) = \text{signe} (\sum_{i=1}^m \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b)$

- Apprentissage :

$$\max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^m \alpha_i$$

$$s.c. \quad \sum_{i=1}^m \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C$$

SVM non linéaire

Formulation duale des SVM

- Modèle : $f(\mathbf{x}) = \text{signe} (\sum_{i=1}^m \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b)$

- Apprentissage :

$$\max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^m \alpha_i$$

$$\text{s.c. } \sum_{i=1}^m \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C$$

Extension non linéaire

- Remplacer $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ par $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j)$

- Modèle : $f(\mathbf{x}) = \text{signe} (\sum_{i=1}^m \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b)$

- Apprentissage :

$$\max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^m \alpha_i$$

$$\text{s.c. } \sum_{i=1}^m \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C$$

Méthodes à noyaux

Espace de Hilbert à noyau reproduisant (RKHS)

- ▶ Espace de Hilbert \approx espace vectoriel \mathcal{H} muni d'un produit scalaire $\langle \cdot, \cdot \rangle_{\mathcal{H}}$
 - ▶ les espaces euclidiens (comme \mathbb{R}^d)
 - ▶ des espaces fonctionnels (comme L_2)

Méthodes à noyaux

Espace de Hilbert à noyau reproduisant (RKHS)

- ▶ Espace de Hilbert \approx espace vectoriel \mathcal{H} muni d'un produit scalaire $\langle \cdot, \cdot \rangle_{\mathcal{H}}$
 - ▶ les espaces euclidiens (comme \mathbb{R}^d)
 - ▶ des espaces fonctionnels (comme L_2)
- ▶ RKHS de noyau reproduisant K : espace de Hilbert \mathcal{H} tel que, pour tout $x \in \mathcal{X}$,
 - ▶ $K(x, \cdot) \in \mathcal{H}$
 - ▶ pour tout $f \in \mathcal{H}$, $\langle f, K(x, \cdot) \rangle_{\mathcal{H}} = f(x)$

Méthodes à noyaux

Espace de Hilbert à noyau reproduisant (RKHS)

- ▶ Espace de Hilbert \approx espace vectoriel \mathcal{H} muni d'un produit scalaire $\langle \cdot, \cdot \rangle_{\mathcal{H}}$
 - ▶ les espaces euclidiens (comme \mathbb{R}^d)
 - ▶ des espaces fonctionnels (comme L_2)
- ▶ RKHS de noyau reproduisant K : espace de Hilbert \mathcal{H} tel que, pour tout $x \in \mathcal{X}$,
 - ▶ $K(x, \cdot) \in \mathcal{H}$
 - ▶ pour tout $f \in \mathcal{H}$, $\langle f, K(x, \cdot) \rangle_{\mathcal{H}} = f(x)$
- ▶ Tout noyau K défini positif engendre un RKHS et on a

$$K(x, x') = \langle K(x, \cdot), K(x', \cdot) \rangle_{\mathcal{H}}$$

$$\mathcal{H} = \left\{ f \in \mathbb{R}^{\mathcal{X}} \mid f = \sum_{i=1}^{\infty} \beta_i K(x_i, \cdot), \beta_i \in \mathbb{R}, x_i \in \mathcal{X}, \|f\|_{\mathcal{H}} < \infty \right\}$$

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \beta_i \beta_j K(x_i, x_j)$$

Méthodes à noyaux

Espace de Hilbert à noyau reproduisant (RKHS)

- ▶ Espace de Hilbert \approx espace vectoriel \mathcal{H} muni d'un produit scalaire $\langle \cdot, \cdot \rangle_{\mathcal{H}}$
 - ▶ les espaces euclidiens (comme \mathbb{R}^d)
 - ▶ des espaces fonctionnels (comme L_2)
- ▶ RKHS de noyau reproduisant K : espace de Hilbert \mathcal{H} tel que, pour tout $x \in \mathcal{X}$,
 - ▶ $K(x, \cdot) \in \mathcal{H}$
 - ▶ pour tout $f \in \mathcal{H}$, $\langle f, K(x, \cdot) \rangle_{\mathcal{H}} = f(x)$
- ▶ Tout noyau K défini positif engendre un RKHS et on a

$$K(x, x') = \langle K(x, \cdot), K(x', \cdot) \rangle_{\mathcal{H}}$$

$$\mathcal{H} = \left\{ f \in \mathbb{R}^{\mathcal{X}} \mid f = \sum_{i=1}^{\infty} \beta_i K(x_i, \cdot), \beta_i \in \mathbb{R}, x_i \in \mathcal{X}, \|f\|_{\mathcal{H}} < \infty \right\}$$

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \beta_i \beta_j K(x_i, x_j)$$

Lien avec l'espace de représentation \mathcal{X}_{ϕ}

- ▶ $\phi : x \mapsto K(x, \cdot)$ et $\mathcal{X}_{\phi} = \mathcal{H}$

Méthodes à noyaux

Formulation fonctionnelle des SVM

- Modèle : $f(x) = \text{signe}(g(x)) = \text{signe}(\langle g, K(x, \cdot) \rangle_{\mathcal{H}})$ avec $g \in \mathcal{H}$
- Apprentissage régularisé par la norme de g :

$$\min_{g \in \mathcal{H}} \frac{1}{2} \|g\|_{\mathcal{H}}^2 + C \sum_{i=1}^m \underbrace{\max \{0, 1 - y_i \langle g, K(x_i, \cdot) \rangle_{\mathcal{H}}\}}_{\text{violation de la marge } = \xi_i}$$

Méthodes à noyaux

Formulation fonctionnelle des SVM

- Modèle : $f(x) = \text{signe}(g(x)) = \text{signe}(\langle g, K(x, \cdot) \rangle_{\mathcal{H}})$ avec $g \in \mathcal{H}$
- Apprentissage régularisé par la norme de g :

$$\min_{g \in \mathcal{H}} \frac{1}{2} \|g\|_{\mathcal{H}}^2 + C \sum_{i=1}^m \underbrace{\max \{0, 1 - y_i \langle g, K(x_i, \cdot) \rangle_{\mathcal{H}}\}}_{\text{violation de la marge } = \xi_i}$$

Théorème de représentation

La solution est une fonction donnée par une combinaison linéaire de noyaux calculés aux points x_i de la base d'apprentissage :

$$g = \sum_{i=1}^m \beta_i K(x_i, \cdot)$$

- Le théorème se généralise à toute fonction de perte et toute régularisation croissante en $\|g\|_{\mathcal{H}}$

Régression non linéaire par Kernel ridge regression (KRR)

Formulation primale dans le RKHS \mathcal{H} (en dimension infinie)

- Modèle $f \in \mathcal{H}$ avec $f(x) = \langle f, K(x, \cdot) \rangle_{\mathcal{H}}$
- Apprentissage :

$$\min_{f \in \mathcal{H}} \sum_{i=1}^m (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

Formulation duale en dimension finie

$$\text{Théorème de représentation} \Rightarrow f = \sum_{i=1}^m \beta_i K(x_i, \cdot)$$

$$\Rightarrow \min_{\beta \in \mathbb{R}^m} \sum_{i=1}^m \left(y_i - \sum_{j=1}^m \beta_j K(x_j, x_i) \right)^2 + \lambda \sum_{i=1}^m \sum_{j=1}^m \beta_i \beta_j K(x_i, x_j)$$

en notation matricielle avec \mathbf{K} la matrice de noyau ($K_{ij} = K(x_i, x_j)$) :

$$\min_{\beta \in \mathbb{R}^m} \|\mathbf{y} - \mathbf{K}\beta\|^2 + \lambda \beta^T \mathbf{K} \beta$$

Solution analytique :

$$\beta = (\mathbf{K}\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}\mathbf{y} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$