

Projet de Master

Romain Mencattini

16 août 2017

Table des matières

1	État de l'art	3
1.1	Introduction	3
1.2	Finance	4
1.2.1	<i>FOREX</i>	4
1.3	Cadre théorique des algorithmes de <i>Machine Learning</i>	6
1.3.1	Introduction	6
1.3.2	<i>Logistic Regression</i>	7
1.3.3	Les arbres de décision	8
1.3.4	<i>Naïve Bayes</i>	9
1.3.5	<i>SVM</i>	10
1.3.6	Réseaux de neurones	12
1.3.7	Descente du gradient	14
1.3.8	<i>Majority vote</i>	17
1.4	<i>Random Subset</i>	18
1.5	<i>Machine Learning</i> dans le cadre de la finance	18
1.6	Conclusion	19
2	Projet	20
3	Bibliographie	21

1 État de l'art

1.1 Introduction

Avant la démocratisation de l'informatique, les opérations financières étaient réalisées par des humains. Ce système pouvait avoir des inconvénients :

- L'émotionnel influençait les transactions. En effet, ces dernières étant effectuées par des humains, il y avait une risque non négligeable que l'état de la personne agisse sur sa décision.
- Un problème sous-jacent était de maintenir une discipline de *trading*. Afin de minimiser les pertes et de maximiser les gains, il fallait se tenir à un plan afin de ne pas se laisser influencer par des paramètres extérieurs. Cela pouvait être très difficile.
- Le *backtesting*¹ était impossible. Tester la qualification ainsi que la qualité de *trading* d'une personne était compliquée. De même pour un *trading plan*.

Ces éléments ont, en partie, favorisé l'émergence et l'utilisation d'algorithmes dans la finance. En 2014 aux USA, 84% des transactions étaient accomplies par des algorithmes [2]. Ce qui représente environ 100000 réalisations, ou *ticks*, par secondes [2]. Durant l'évolution l'outil informatique, le monde de la finance en a suivi les améliorations afin de perfectionner leurs algorithmes. On retrouve donc des méthodes d'optimisations poussées ainsi que les récentes découvertes de *data mining* et de *machine learning*, abrégé *ML*. Des propositions de plus en plus pointues dans les deux domaines voient le jour. L'algorithme qui sera au coeur de ce projet en fait partie. Il s'agit d'un réseau de neurones avec plusieurs couches prenant en compte des paramètres particuliers à la finance.

Afin d'approcher aux mieux ces notions, nous allons discuter des éléments nécessaires à leur compréhension. Nous allons en premier lieu traiter une partie domaine financier ainsi que ces outils. Puis nous parlerons de plusieurs méthodes de *ML*. Voici celles vont être développées dans cet état de l'art :

- Les réseaux de neurones.
- Les arbres de décision.
- Les algorithmes *SVM* [7].
- *Logistic Regression*.
- *Naive Bayes*.
- Descente du Gradient [9] ainsi que sa version dite stochastique [8].

Finalement, nous lierons les deux domaines en montrant comment adapter les modèles mathématiques de *ML* pour les utiliser comme techniques de *trading*, en évaluant leur performances.

1. Backtesting is the process of testing a trading strategy on relevant historical data to ensure its viability before the trader risks any actual capital. [15]

1.2 Finance

1.2.1 FOREX

Afin d'appréhender le fonctionnement du *FOREX*, il est important de mentionner certaines décisions historiques. Ces dernières ayant façonné le marché des devises actuel.

Jusqu'à la première guerre mondiale, le système en vigueur se basait sur l'or, que l'on nommait l'étalon-or¹. S'en suit une période d'instabilité notamment due aux pertes occasionnées par la guerre, un après-guerre compliqué, la crise boursière de 1929 et la seconde guerre mondiale.

C'est au sortir de cette dernière, que la nécessité de "*mettre en place une organisation monétaire mondiale et de favoriser la reconstruction et le développement économique des pays touchés par la guerre*" [17], est apparue. Le but était également "*d'aplanir les conflits économiques, reconnaissant par là les problèmes engendrés par les disparités économiques*" [18].

Plusieurs idées furent proposées, mais ce fût celle de Harry Dexter White qui fût mise en place. Cette dernière prévoyait entre autre :

- le choix du Dollar américain comme étalon, avec rattachement à l'or².
- Création de la Banque internationale pour la reconstruction et le développement (BIRD) qui deviendra la banque mondiale.
- Le Fond monétaire international (FMI).
- Création de l'Organisation mondial du commerce³.

On remarque que ces institutions sont toujours en activités, cela démontre l'importance de ces accords pour le système financier actuel. Cela est également vrai pour le marché des taux de changes, où le USD est toujours utilisé entre deux échanges. Il n'est en effet pas possible de faire CHF/EUR⁴. Le passage par le USD entre les deux est obligatoire ; nous aurons donc CHF/USD⁵ puis USD/EUR⁶, même si cela reste transparent pour l'utilisateur.

Le marché *FOREX* porte sur les devises. La valeur d'une devise ne peut être exprimée qu'en fonction d'une autre. Par exemple 1 franc suisse vaut 1.05 euro.⁷ La transaction porte donc sur deux monnaies comme CHF/EUR. On va vendre des francs suisses pour acheter euros ou l'inverse. Le nom du marché vient d'ailleurs de ces échanges. On échange une monnaie contre une autre, c'est un *FOreign EXchange*, ou *FOREX*.

Il y a deux variations possibles :

- La monnaie peut subir une dépréciation.

1. Source : [18].

2. Suspension l'équivalence or pour le dollar américain en août 1971 puis abandon définitif en mars 1973 [17].

3. Ne verra le jour qu'en 1995 faute d'accord [17].

4. Franc Suisse - Euro

5. Franc Suisse - Dollar Américain

6. Dollar Américain - Euro

7. Taux fictif utilisé pour l'exemple.

— La monnaie peut subir une appréciation.

Lorsque le prix d'une devise augmente par rapport à une monnaie étrangère, on parle d'appréciation. Ainsi dans le cas contraire, on parlera d'une dépréciation.

La mondialisation a facilité ce marché. En effet, toutes devises étant accessibles depuis n'importe où, il devient donc possible d'avoir des marchés avec des devises plus exotiques.

Les principaux acteurs financiers sont [19] :

- Les banques commerciales. Elles peuvent pratiquer des interventions directes car elles gèrent des dépôts et veulent opérer des transactions sur ces derniers. Il leur est également possible de réaliser le rôle d'intermédiaire financier.
- Les entreprises. Ces dernières vont pratiquer des transactions directes, si elles disposent d'un accès aux marchés sinon via des intermédiaires.
- Les institutions financières non-bancaires. On peut citer les fonds de pensions, les sociétés d'assurances ou les *hedge funds*. Ce sont surtout dans un but de spéculation, d'arbitrage ou de couverture de risque qu'elles agissent.
- Les banques centrale. Il peut y avoir des interventions directes, dans le but de modifier l'appréciation de la monnaie.
- Les ménages. Surtout dans une optique de voyage, d'achat ou de spéculation.

Henry Bourguinat a énoncé "*la règle des trois unités*" qui correspondent aux unités de temps, de lieu et d'opérations et d'acteurs. Le *FOREX* répond à ces trois unités [20] :

- Ce marché fonctionne 24h/24 et les transactions s'effectuent presque en continue.
- Il fonctionne à l'échelle mondiale tout en étant décentralisé. De part l'évolution des technologies, l'information circule aisément malgré son statut.
- L'uniformité des procédés ainsi que des produits est présente. Les acteurs malgré nationalité sont de même nature.

Il existe principalement deux horizon temporels : le *spot* et le *forward*.

Le premier est également appelé "Le marché au comptant". Lorsque deux acteurs se mettent d'accord sur une transaction, cette dernière se réalise immédiatement ¹.

Le second peut être nommé "Le marché à terme". L'accord est passé à un temps T mais la transaction effective ne se réalise que dans le futur. Ce futur, ou maturité, peut être de plusieurs dizaines de jours, voir des années, soit $T + X$ ².

Il y a opérations réalisable sur le marché à terme. :

- Les *swaps*. Ils consistent à vendre une monnaie au comptant puis à la racheter à terme ³.
- Les *futures/forwards*. La différence entre ces deux tient surtout à leur standardisation et leur mise en place. Cependant le principe reste le même : on réalise une opération (d'achat ou de vente) qui ne s'effectuera qu'à maturité.

1. Valable en théorie, dans la réalité cela peut prendre du temps [19]

2. Où T est le moment présent, et X une durée de temps.

3. Soit à $T + X$

- Les *options*. Cela représente un contract vendu par un parti (*the option writer*) à un autre parti (*the option holder*). Ce contrat offre le droit, et non l'obligation contrairement aux *futures/forwards*, d'acheter (*call*) ou de vendre (*put*). Ici encore, il faut attendre la maturité.

Les options sont très versatiles. Elles peuvent être utilisées afin de spéculer ou de diminuer le risque. Voici les différents types possibles :

- **long call** → on achète le droit d'acheter le sous-jacent à un certain prix.
- **short call** → on vend le droit d'acheter le sous-jacent à un certain prix.
- **long put** → on achète le droit de vendre le sous-jacent à un certain prix.
- **short put** → on vend le droit de vendre le sous-jacent à un certain prix.

Le *bid* est le prix maximum qu'un acheteur est d'accord de payer pour un sous-jacent. De la même manière, le *ask* est le prix minimum qu'un vendeur accepte pour vendre un sous-jacent [21]. Lorsqu'il y a un *bid* qui a la même valeur qu'un *ask*, une vente est effectuée.

La différence entre le *bid* et le *ask*, appelée le *spread*, représente la liquidité d'un actif. Il est également utilisé comme marge par les *broker*[22] et autres plateformes.

1.3 Cadre théorique des algorithmes de *Machine Learning*

1.3.1 Introduction

T. Mitchell a donné une définition formelle [6] :

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P is its performance at tasks in T , as measured by P , improves with experience E "

On a donc une tâche T à accomplir, où T peut consister trier des images ou à reconnaître des motifs. La mesure de la réussite de cette tâche T est nommée P . C'est-à-dire la qualité du résultat du programme pour la tâche donnée, T . Si le programme améliore son résultat P pour la tâche T grâce à de l'expérience E . Il s'agit d'un programme de *machine learning*. L'expérience peut être vu comme une phase d'entraînement ou comme le fait de retenir les réponses après avoir accompli la tâche.

Il existe deux catégories d'apprentissage :

- L'apprentissage supervisé.
- L'apprentissage non-supervisé.

Dans le cas du premier, on fournit au programme, un ensemble d'entraînement¹, qui contient des réalisations ainsi que le résultat de la classification. Le programme va donc pouvoir utiliser ce savoir afin d'améliorer sa performance P . Nous disposons donc de nombreux couples (x_i, y_i) et le but est de trouver une fonction $f \in F$ telle que : $f(x) = y$.

Pour l'apprentissage non-supervisé, on fournit des données, mais sans le résultat voulu. C'est uniquement après avoir décidé d'une valeur qu'on va signifier au programme si cette

1. Ou d'expérience, E

dernière est correcte. On ne lui donnera jamais la valeur attendue. Il va donc utiliser uniquement les résultats précédents pour améliorer son P .

Par exemple, on désire reconnaître un certain type de voiture à partir d'images. Dans le cas de l'apprentissage supervisé, nous allons fournir au programme un ensemble d'entraînement qui contient de nombreuses photos de voitures, ainsi que la marque des dites voitures. L'algorithme va donc travailler avec ces données.

Par contre dans le cas de l'apprentissage non-supervisé, le programme ne pourra utiliser que les photos, et après avoir retourné le résultat, nous lui dirons si c'est juste ou faux. Il mémorisera le résultat et va tenter d'optimiser ses réponses.

Concernant, l'ensemble d'entraînement, il y a des points à prendre en compte afin de minimiser les risques de sur-apprentissage¹, et de maximiser la qualité de nos données. Pour ce faire il faut :

- Représenter la population générale. Donc si le but est du traitement de la langue, il faut que la propension et la répartition des mots soient les mêmes que ceux de la langue.
- Contenir des membres de chaque classes. Pour reconnaître des chiffres, il est important de disposer de chacun des chiffres dans l'ensemble d'entraînement.
- Contenir de grandes variations ainsi que du bruit. Afin d'éviter le sur-apprentissage, il faut de nombreux exemples différents, voir très différents, les uns des autres ainsi que du bruit².

Il est important de saisir comment fonctionne les algorithmes de *machine learning*. Les données étant "fixes", elles ne peuvent être modifiées afin d'augmenter la performance de classification. Une fois qu'on dispose d'un *tick* en finance ou bien d'une *review* sur un produit, on ne peut en changer l'expressivité. Il convient donc d'utiliser ces données, souvent de très hautes dimensionnalités³, dans des équations dont on pourra faire varier les paramètres afin de classifier au mieux. Le coeur des algorithmes de *machine learning* se trouve donc dans ses équations, qu'il va convenir d'optimiser.

1.3.2 Logistic Regression

Une régression en statistique consiste à analyser la relation entre une variable par rapport à un ensemble d'autres [23]. On veut estimer la probabilité conditionnelle, en se basant sur des variables et en utilisant une distribution logistique cumulative. Cette dernière a une forme semblable à une distribution Gaussienne, mais avec des queues épaisses et donc une *kurtosis* plus élevée. La *kurtosis* étant définie comme suit :

$$Kurt(X) = \frac{\mu_4}{\sigma^4}, \text{ où } \mu_4 \text{ est le quatrième moment centré et } \sigma^4 \text{ la variance au carré.}$$

1. Le sur-apprentissage consiste à apprendre par coeur la tâche, plutôt que d'apprendre les principes pour réaliser la tâche.

2. Comme des faux exemples.

3. Cela signifie qu'elles sont représentées par un grand nombre d'attributs.

Le but est de modéliser : $P(Y = 1|X = x)$ comme fonction de x . Nous voulons donc savoir quelle est la probabilité que la classe Y vaille 1 sachant que X vaut x .

Le modèle de régression est le suivant [3] :

$$\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + x \cdot \beta$$

En résolvant cette équation pour p , cela donne [3] :

$$p(x|y) = \frac{1}{1+e^{-(\beta_0+x\cdot\beta)}}$$

Dans le cadre de l'article : *A Machine Learning Approach to Automated Trading* [3], l'auteur a implémenté deux variations de cet algorithmes :

— *Logistic regression with a ridge penalty* :

$$\sum_{i=1}^N (y_i - \sum_j \beta_j x_{ij})^2 + \lambda \sum_j \beta_j^2$$

L'objectif est de minimiser le carré de la différence entre la classe y_i et le résultat calculé : $\sum_j \beta_j x_{ij}$. Donc en fonction de l'observation x_i et du coefficient β . En

ajoutant une pénalité, sur β , on va tenter d'éviter le sur-apprentissage.

— *Lasso logistic regression/ Lasso regularization* :

$$\sum_{i=1}^N (y_i - \sum_j \beta_j x_{ij})^2 + \lambda \sum_j |\beta_j|$$

Le fonctionnement est similaire au précédent, seul change la pénalité.

Comme mentionné plus haut, l'optimisation concerne les paramètres variables. Dans le cas de la régression logistique, il s'agit du β . Il conviendra donc de trouver la valeur optimale pour cette variable afin de maximiser ou minimiser les équations ci-dessus. Cette manière de faire s'applique aux autres techniques de *ML*. Les équations et les paramètres vont changer, mais le but sera toujours d'optimiser ces éléments.

1.3.3 Les arbres de décision

Un arbre de décision est un arbre, dont chaque nœud représente un test sur un attribut. Les branches qui suivent directement le nœud sont les valeurs possibles de l'attribut. Les feuilles de l'arbre, quant à elles, sont la classification de l'élément donné en entrée.

Il est important de disposer des attributs avant de commencer la construction de l'arbre. Lors de l'implémentation, nous pouvons représenter l'arbre comme une suite de *if-then-else* afin d'améliorer la lisibilité. Dans ce cas très précis, disposer d'un langage permettant le *pattern matching* est fort utile.

Cet algorithme a tendance à très facilement sur-apprendre, il convient donc de bien choisir la manière de construire l'arbre ainsi que l'ensemble d'entraînement pour minimiser cet effet.

Voici un exemple d'arbre de décision¹ :

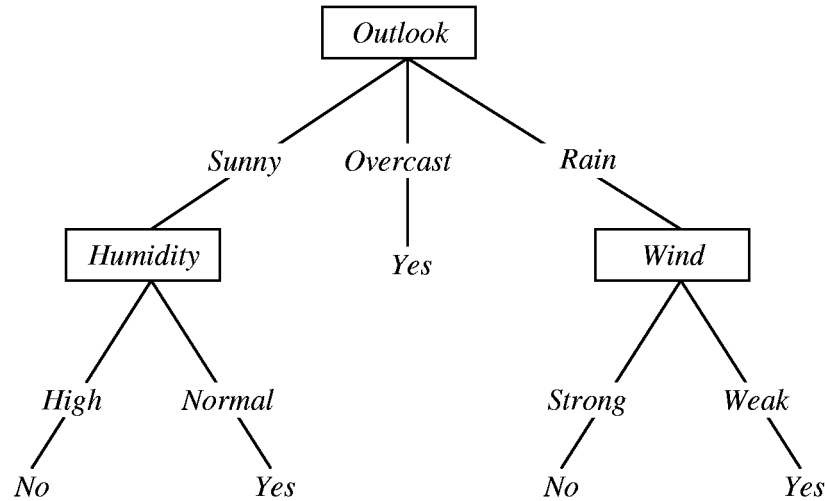


FIGURE 1 – Exemple d'arbre de décision : Permet de décider si nous pouvons aller jouer au tennis ou non.

Afin de construire l'arbre à partir de l'ensemble d'entraînement, il existe plusieurs algorithmes. Un des plus connus est le *ID3* [10]. Il s'agit d'une méthode de type *greedy* qui va à chaque itération, effectuer un test statistique² afin de savoir quel attribut est le plus discriminant, utiliser cet attribut comme nœud, puis itérer jusqu'à avoir utilisé tous les attributs.

1.3.4 Naïve Bayes

À l'instar de la régression logistique [1.3.2], il s'agit d'un classifieur probabiliste. Ce dernier se base sur le théorème de Bayes³ :

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

Il est important que les données aient une distribution normale.

À partir de cela, nous obtenons l'équation de *machine learning* suivante [3] :

$$V_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i|v_j)$$

1. <http://cloudmark.github.io/images/kotlin/ID3.png>

2. Il vise à vérifier la quantité d'information gagnée pour la classification [10].

3. <https://brilliant.org/wiki/bayes-theorem/>

Où V_{NB} est la classe obtenue, $P(v_j)$ la probabilité à *priori* donc sans informations et $\prod_i P(a_i|v_j)$ la probabilité de vraisemblance. Il convient donc de trouver la classe qui maximise ce calcul.

Afin d'avoir une certaine sécurité dans les résultats, il est possible d'ajouter un seuil. Les réponses du classifieur étant comprise entre 0 et 1, le seuil permettra de décider si la réponse sera prise en compte.

Par exemple, avec un seuil de 0.6 si $V_{NB} = 0.58$, alors la réponse n'est pas validée. Si $V_{NB} = 0.7$ la réponse est jugée sûre.

La *ROC¹ Curve analysis* permet d'améliorer le classifieur. En effet cette dernière peut détecter les *true positive rate* par rapport aux *false positive rate* pour différents seuils de classification de l'algorithme *Naive Bayes* [3]. À partir de cela, nous pouvons déterminer le meilleur seuil de sortie et donc perfectionner notre algorithme. De plus cette courbe peut aider à comparer des classifieurs entre eux en comparant la surface sous la courbe [3].

La méthode utilisée est la suivante [3] : il faut faire s'intersecter la pente S avec la courbe *ROC* et ainsi obtenir une valeur optimale pour le seuil. Cette pente S est définie comme suit [3] :

$$S = \frac{Cost(P|N) - Cost(N|N)}{Cost(N|P) - Cost(P|P)} \cdot \frac{N}{P}$$

Sachant que $Cost(P|N)$ est le coût pour avoir mal classé une classe négative comme positive. P est la somme des vrais positifs et des faux négatifs. N , quant à lui, vaut la somme des vrais négatifs ainsi que des faux positifs.

Voici une illustration² :

1.3.5 SVM

Le but de l'algorithme *SVM*³ est de séparer les données grâce à un hyper-plan. Cela permet de différencier les classes des observations suivantes en déterminant s'ils se trouvent d'un côté ou l'autre du plan séparant. Ce plan n'est pas unique⁴ :

Notre fonction est :

$$f(x) = (w \cdot x) + b$$

Le but est de maximiser la distance entre les points les plus proches de l'hyper-plan, tout en pénalisant les points mal classés. Il n'est pas toujours possible de séparer les données de dimensions n , il conviendra donc d'augmenter la dimension afin d'obtenir une dimension $m > n$ plus discriminante. La fonction $\phi(x)$ est utilisée dans ce but. Un exemple de fonction est :

1. Receiver Operating Characteristic

2. Source : http://www.prolekare.cz/dbpic/jp_5403_f_20-x1000_1600

3. *Support Vector Machine*

4. Source : <https://computersciencesource.files.wordpress.com/2010/01/svmafter.png>

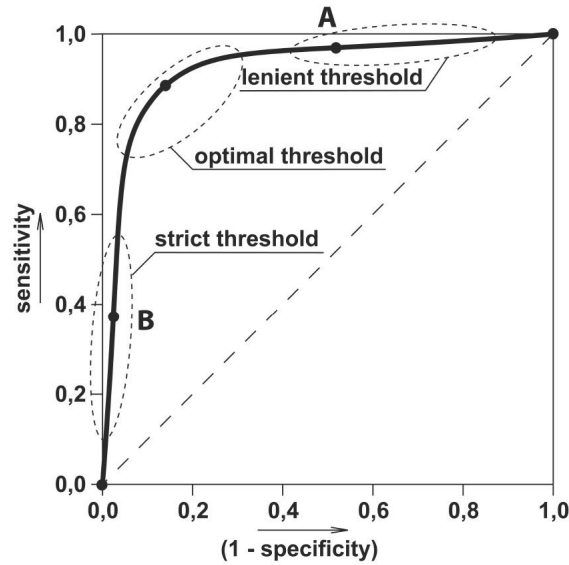


FIGURE 2 – Exemple de *ROC curve* : L'axe des x correspond au taux de faux positifs et l'axe des y au taux de vrais positifs. On veut donc tendre vers $y = 1$ et $x = 0$. Lors de l'optimisation par la pente S , le but sera d'obtenir une intersection avec la *ROC curve* dans la zone *optimal threshold* afin d'avoir le meilleur seuil possible.

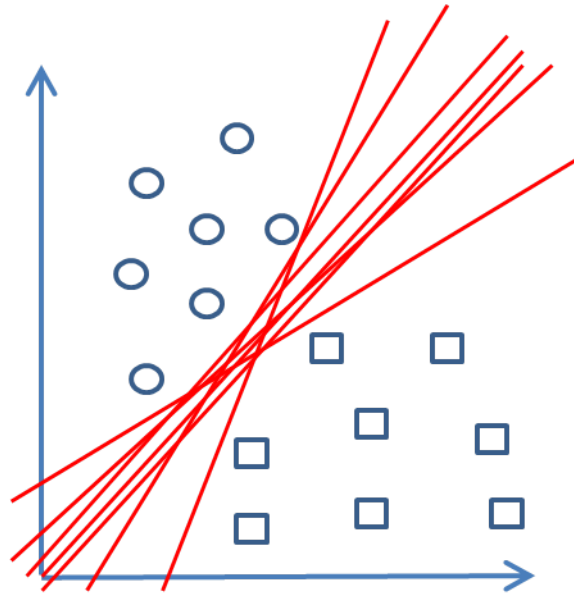


FIGURE 3 – Exemple d'hyper-planes séparant les données. Si nous voulons classifier un nouvel élément, il suffit de calculer s'il se trouve à gauche ou à droite de l'hyper-plan. Dans le premier cas, il s'agira, pour notre algorithme, d'un rond et dans l'autre d'un carré

$$\phi : R^2 \rightarrow R^3(x, y) \rightarrow (x, y, z) := (x^2, \sqrt{2}xy, y^2)$$

Il est possible d'imager cette opération comme cela ¹ :

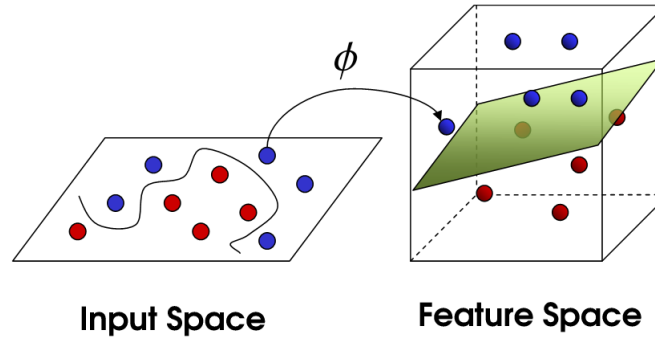


FIGURE 4 – Exemple d'utilisation de la fonction ϕ pour passer d'un espace R^2 à R^3 afin de faciliter la séparation.

En terme d'équation, nous voulons minimiser w , soit $\|w\|^2$ dans l'équation de l'hyperplan : $(w \cdot x) + b$. Ce qui donne :

$$y_i(w \cdot x_i + b) \geq 1 \text{ si } y \in \{-1, +1\} \text{ [12]}$$

Si malgré l'augmentation de la dimension, les données ne sont pas séparables, il faut tenter de minimiser le nombre d'éléments mal placés. Pour ce faire :

$$y_i(w \cdot x_i + b) \geq 1 - \xi_i \text{ avec } \xi_i > 0 \text{ [12]}$$

Afin de diminuer l'erreur et optimiser au mieux notre classifieur.

1.3.6 Réseaux de neurones

Tout comme les algorithmes génétiques s'inspirent de la sélection naturelle dans un but d'optimisation, les réseaux de neurones se basent sur un modèle formels de neurones ² afin de copier la capacité d'apprentissage des êtres vivants.

Il s'agit d'opérer à partir de données en entrée, des *inputs*, une ou plusieurs multiplications matricielles en utilisant des vecteurs de poids, des *weights*. L'optimisation s'applique sur les *weights*, afin de maximiser la classification.

Un réseau de neurones peut avoir plusieurs couches, *layers*. Dans ce cas, la première couche est appelée *input layer*, la dernière *output layer* et toutes celles entre ces deux sont les *hidden layers*. De plus les neurones peuvent être pleinement connectés avec ceux de la couche suivante, *feed-forward* ; ce qui signifie que les neurones de la couche $n-1$ influencent

1. <https://www.dtreg.com/uploaded/pageimg/SvmDimensionMap.jpg>

2. Neurones formels : <http://www.peoi.org/Courses/Coursesfr/neural/neural3.html>

ceux de la couche n . Il est également possible que les éléments de n agissent sur les neurones de $n - 1$, ce phénomène est appelé *feedback networks*.

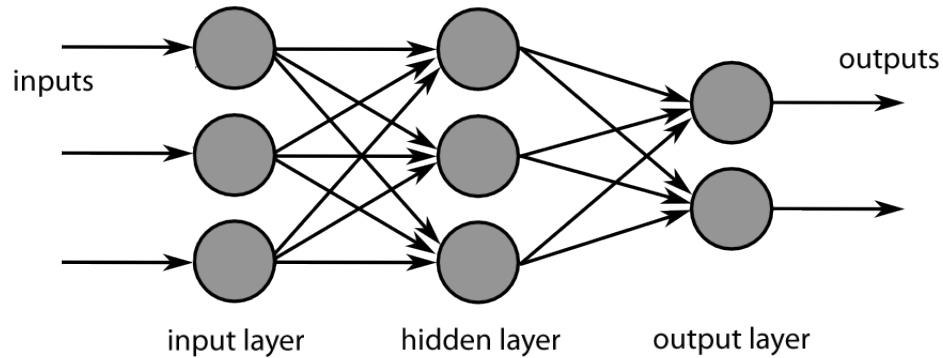


FIGURE 5 – Exemple de réseau de neurones avec plusieurs couches. Illustre également le *feed-forward* : l'output de l'*input layer* est propagé dans chaque neurones de l'*hidden layer*. Même chose pour les deux dernières couches.¹

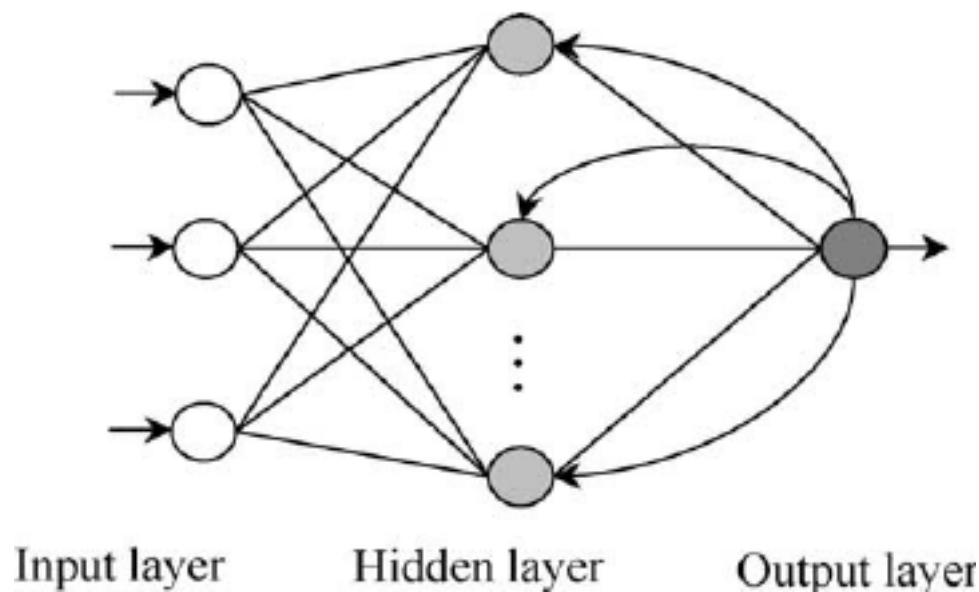


FIGURE 6 – Exemple de réseau de neurones avec plusieurs couches. Dans ce cas là, les couches font du *feed-forward* mais de plus, le *feedback* est utilisé afin d'influencer les couches précédant l'*output layer*.³

1. Source : http://web.utk.edu/~wfengl/spark/_images/fnn.png

3. Source : <https://www.researchgate.net/profile/Yen-Ming-Chiang/publication/222562767/figure/fig2/AS:305112384851969@1449755869325/Fig-2-The-architecture-of-a-recurrent-neural-network.png>

Concernant la valeur de sortie, elle peut être simple, comme le résultat d'une classification binaire, *i.e.* 0 ou 1 en sortie. Mais elle peut également être d'une dimensionnalité plus élevée comme un vecteur, citons l'exemple d'un point dans un espace \mathbb{R}^2 .

Afin de borner les valeurs en sortie, la plupart des réseaux utilisent une fonction. Nous pouvons citer :

- La fonction sigmoïde : $S(t) = \frac{1}{1+e^{-t}}$.
- La fonction tangente hyperbolique : $f(x) = \tanh(x)$.

Bien souvent, l'utilisation d'un réseau de neurones à une couche est suffisante. Cela est valable pour les fonctions continues, dans le cas de fonctions discontinues, il est intéressant de passer à un réseau disposant de plusieurs couches. Attention toutefois, si le nombre de neurones est trop important, l'algorithme va avoir tendance à sur-apprendre, et à l'inverse, à sous-apprendre si le nombre est trop faible. Il est donc important de bien doser cette quantité afin d'éviter ces problèmes.

L'article qui se trouve au coeur du projet [1] propose un algorithme basé sur les réseaux de neurones. Il s'agit de trois couches qui doivent optimiser chacune une série de paramètres. La première va maximiser la première série de paramètres suivant certains critères. Une fois ces paramètres changés, ils deviennent des constantes pour les deux prochaines couches. La deuxième va améliorer les paramètres qui lui sont attachés en gardant les valeurs des paramètres de la première couche constantes. Le même principe est répété pour la troisième couche.

Ce n'est donc pas un réseau à multiples couches, dans le sens où, une série de paramètres n'est optimisée que par une couche et non plusieurs. On peut visualiser cela comme trois réseaux qui chaînent leurs résultats.

1.3.7 Descente du gradient

L'algorithme de descente du gradient fonctionne sur des fonctions réelles différentiables sur un espace tel que \mathbb{R}^n . Il est itératif et fonctionne donc en améliorant l'itération précédente, jusqu'à atteindre une condition d'arrêt.

Avant d'en expliquer la teneur mathématique, voici un exemple de l'algorithme dans \mathbb{R}^2 :

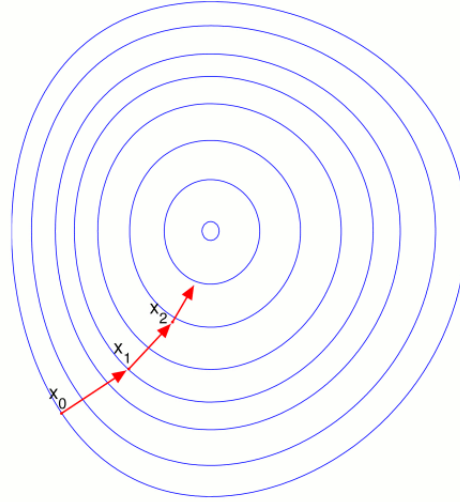


FIGURE 7 – Exemple de l’algorithme de descente du gradient en deux dimensions. À chaque itération, il faut prendre la direction opposée au gradient¹, cela permet d’arriver à un nouveau point. En itérant, on se rapproche de l’maximum.

Afin de comprendre, les divers algorithmes, il est important d’avoir les connaissances mathématiques sur ce sujet. L’algorithme se définit comme suit² :

Soit un point initial $x_0 \in \mathbb{R}$. Soit $\epsilon > 0$ un seuil de tolérance. L’algorithme définit une suite d’itération $x_1, x_2, \dots \in \mathbb{R}^n$, jusqu’à ce qu’un test d’arrêt soit satisfait. Pour passer de x_i à x_{i+1} , il faut :

- Calculer $\nabla f(x_k)$
- Si $\|\nabla f(x_k)\| \leq \epsilon$ alors arrêt.
- Sinon il faut calculer α_k par recherche linéaire sur f en x_k . Cette recherche se fait dans la direction opposée au gradient, soit $-\nabla f(x_k)$. Une fois α_k calculé, il faut mettre à jour le point itéré :

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

La preuve que la recherche dans la direction opposée au gradient induit une décroissance est la suivante. Si la dérivée est non nulle³ au point x , $f'(x) \neq 0$. Soit

$$d = -\nabla f(x)$$

Puisque :

$$f'(x) \cdot d = \nabla f(x) \cdot -\nabla f(x) = -\|\nabla f(x)\|^2 < 0$$

1. Dans cet exemple, il est possible de dire, qu’il faut prendre la normal de la courbe de niveau.

2. Inspiré de [9]

3. i.e : Si nous ne sommes pas déjà sur un maximum

L'égalité est strictement petit car la dérivée est non nulle par hypothèse. Cela implique que :

$$f(x - \alpha \nabla f(x)) < f(x), \forall \alpha > 0$$

Nous avons donc que pour chaque itération, la valeur obtenue va décroître jusqu'à atteindre un maximum ou le seuil ϵ .

La descente du gradient d'un point de vue mathématique peut également être utilisée conjointement à un réseau de neurones. En effet dans cet article [1], les poids w de la première couche, sont optimisés suivant cette formule :

$$w_{i,t} = w_{i-1,t} + \rho \Delta w_{i,t}$$

Où t correspond au temps, $i - 1$ à l'itération courante et ρ le taux d'apprentissage. De part la convergence de w_i vers son optimale, cela permet, une fois cette valeur obtenue, de l'injecter comme étant les poids d'un réseau de neurones.

Il existe deux types d'algorithme de descente du gradient :

- Le *Batch Algorithm* a pour but de minimiser la fonction de coût qui aura été définie au préalable. Ce dernier peut prendre énormément de temps sur des gros jeux de données de par son côté itératif.
- Le *Online Algorithm* va prendre les exemples un à un. L'optimisation ne va se faire que pour cet unique exemple puis va recommencer sur une autre donnée. Cette manière de procéder s'applique aisément sur les gros volumes de données.

Tiré cet article [8], voici des équations pouvant être minimisée dans le cadre d'un programme de *machine learning*. Soit (x_i, y_i) , $i = 1..N$, un ensemble de données. Notre fonction de coût est $l(y, y')$. Cela représente le coût de prédire y' quand la réponse est y . Moyennons cela sur tous les exemples :

$$E_N(f) = \frac{1}{N} \sum_{i=1}^n l(f_w(x_i), y_i)$$

Où f_w est une fonction pondérée par un vecteur de poids w que l'on cherche à optimiser. Pour optimiser les poids, nous allons utiliser la descente du gradient :

$$w_{k+1} = w_k - \alpha_k \nabla f_{w_k}(x) \rightarrow w_{k+1} = w_k - \alpha_k \sum_{i=1}^n \nabla Q((x_i, y_i), w_k)$$

Où $Q(x, y) = l(f_w(x), y)$.

Comme mentionné plus haut, la technique de *Batch Algorithm* devient lente lorsque le nombre données augmente. Pour pallier à ce problème, il existe une technique : *stochastic gradient descent*. C'est une méthode d'approximation statistique de la descente du gradient.

Il convient de ne prendre qu'un seul couple au lieu de l'ensemble complet d'entraînement. La somme disparaît donc dans l'équation à minimiser :

$$w_{k+1} = w_k - \alpha_k \nabla Q((x_i, y_i), w_k)$$

C'est donc une technique de *Online Algorithm*. De part le fait qu'il n'a pas besoin de se souvenir des exemples précédant, l'algorithme peut traiter des données à la volée.

Il est possible pour pénaliser la complexité de w de rajouter un élément. Cela d'essayer de borner les valeurs des poids :

$$w_{k+1} = w_k - \alpha_k \nabla Q((x_i, y_i), w_k) + \sigma P(w_k)$$

Où $\sigma > 0$ est un hyper-paramètre et P peut valoir :

- **L1 norm** = $P(w) := \sum_{i=1}^n |w_i|$
- **L2 norm** = $P(w) := \frac{1}{2} \sum_{i=1}^n w_i^2$
- **Elastic Net** = $P(w) := \rho \frac{1}{2} \sum_{i=1}^n w_i^2 + (1 - \rho) \sum_{i=1}^n |w_i|$

1.3.8 Majority vote

Le *majority vote* n'est pas une algorithme de *ML* à proprement parlé. Il s'agit plutôt d'améliorer les résultats en modifiant la manière d'utiliser les algorithmes pré-mentionnés. Il est donc possible d'utiliser le *majority vote* avec la régression logistique, les réseaux de neurones, etc.

Le principe de l'algorithme est le suivant. Soit E notre ensemble d'entraînement, soit $M = \{f_1, f_2, \dots\}$ notre ensemble de méthodes f_i de *machine learning*. Notons $M_E = \{f_{1_E}, f_{2_E}, \dots\}$ notre ensemble de méthodes f_{i_E} entraînées.

$$f_{i_E}(x) = \begin{cases} 1 & \text{si la classe calculée de } x \text{ est } 1. \\ 0 & \text{si la classe calculée de } x \text{ est } 0. \\ -1 & \text{si aucune classe n'est trouvée ou si le seuil est insuffisant.} \end{cases}$$

L'algorithme se comporte comme suit pour une observation x :

$\forall f_{i_E} \in M_E$, il faut calculer $f_{i_E}(x)$ et garder la classe résultante dans notre liste de résultat R . Si $\exists classe_i \in R$, telle que

$$\sum (classe_i \in R) > \frac{|M_E|}{2}$$

Alors cela signifie que la $classe_i$ est notre résultat, dans le cas contraire, l'algorithme ne donne aucune réponse.

L'algorithme va donc, à partir d'un ensemble de méthodes et d'une observation, calculer les classes. Si une de ses dernières est représentée de manière majoritaire¹, le programme retournera ce résultat. Si la majorité n'est pas atteinte, il n'a aucune classe jugée valable et donc aucun résultat ne sera rendu.

Cette technique réduit la composante individuelle des méthodes de *ML* ainsi que le risque d'erreur. Si une observation est mal classifiée, cela signifie que plus de la moitié des algorithmes ont fait une erreur. Dans ce cas là, il convient de changer les attributs utilisés ou l'ensemble d'entraînement car l'erreur a peu de chance de venir d'un problème d'implémentation ou d'un algorithme particulier. Hors les seules composantes communes sont l'ensemble d'entraînement ainsi que les attributs.

1.4 *Random Subset*

À l'instar de la technique de *majority vote*, le *random subset* est un algorithme pour utiliser des méthodes de *ML*.

Le principe est le suivant :

Il faut prendre N différents ensembles de données² du domaine concerné. Ces N ensembles seront nos ensembles d'entraînements. Ils permettront d'entraîner un algorithme dessus³. Nous disposerons donc de N instances de l'algorithme de *ML* initial, mais avec un entraînement différent pour chacun, puis pour chaque observation, nous allons utiliser ces N instances afin d'obtenir un résultat de classification.

Comme pour le *majority vote*, si une même classe est représentée une majorité de fois, le *random subset* retournera ce résultat, et aucun dans le cas contraire.

Il est nécessaire d'avoir un grand jeu de données afin de fournir suffisamment d'échantillons aux instances pour qu'elles apprennent correctement. Une fois ce problème résolu, cette technique nous donne la possibilité d'utiliser au mieux l'ensemble d'entraînement. En le fractionnant, cela permet d'entraîner les algorithmes avec des données différentes et donc d'augmenter l'horizon de connaissance du programme.

De plus, le système de majorité promet une qualité et une confiance accrue dans les résultats obtenus.

1.5 *Machine Learning* dans le cadre de la finance

Avant de parler des performances des algorithmes mentionnés, il convient de préciser qu'ils proviennent de plusieurs articles et recherches différents. Ils ont donc été testés avec des paramètres ayant des valeurs disparates ainsi que sur des données distinctes. Il est donc très compliqué de comparer les résultats des algorithmes et d'affirmer qu'une méthode est meilleure qu'une autre. Ces résultats servent plutôt à illustrer la qualité intrinsèque des

1. *i.e.* 50% + 1 voix

2. Il est possible de découper notre ensemble initial afin d'atteindre ce critère.

3. Comme *Naïve Bayes* ou SVM.

procédés. Avec un résultat de 80%, nous pouvons estimer que la performance est bonne, et l'inverse pour une valeur de 20%.

1.6 Conclusion

2 Projet

3 Bibliographie

1. "An automated fx trading system using adaptive reinforcement learning" de M.A.H Dempster et V. Leemans 2004.
2. Financial Times : "Real investors eclipsed by fast trading", 2012 <https://www.ft.com/content/da5d033c-8e1c-11e1-bf8f-00144feab49a?mhq5j=e1>
3. "A Machine Learning Approach to Automated Trading", 09.05.2016, Ning Lu
4. "An efficient implementation of the backtesting of trading strategies." Ni, Jiarui, et Chegqi Zhang, *Parallel and Distributed Processing and Applications* (2005) : 126-131.
5. "Algorithmic Trading : Winning Strategies and Their Rationale (Wiley Trading Series)", John Wiley and Sons, 2013
6. "Machine Learning", Mitchell, Tom M. New York, 1997.
7. Article Wikipédia sur SVM : https://fr.wikipedia.org/wiki/Machine_%C3%A0_vecteurs_de_support
8. "Online Machine Learning Algorithms For Currency Exchange Prediction", Eleftherios Soulas et Dennis Shasha de NYU, Courant Department.
9. Article Wikipédia sur Algorithme du gradient : https://fr.wikipedia.org/wiki/Algorithme_du_gradient
10. "Descision Tree Learning", Tom M. Mitchell
11. Article Investopedia sur les Options <http://www.investopedia.com/terms/o/option.asp>
12. "Support Vector Machine (and Statistical Learning Theory) Tutorial", de Jason Weston, NEC Labs America. http://www.cs.columbia.edu/~kathy/cs4701/documents/jason_svm_tutorial.pdf
13. "An Introduction to Neural Networks" Vincent Cheung et Kevin Cannons : <http://www2.econ.iastate.edu/tesfatsi/NeuralNetworks.CheungCannonNotes.pdf>
14. Exemple de réseaux de neurones <http://csc.lsu.edu/~jianhua/nn.pdf> p.5
15. *Backtesting* Investopedia <http://www.investopedia.com/terms/b/backtesting.asp>
16. Historique des taux de changes : <http://www.xe.com/currencycharts/>
17. Article Wikipedia sur les Accords de Bretten Woods : https://fr.wikipedia.org/wiki/Accords_de_Bretton_Woods
18. Article Étalon-Or à Étalon-Dollar : <http://la-chronique-agora.com/etalon-or-etalon-dollar/>
19. *CHAPITRE 1 : LE MARCHE DES CHANGES Monnaie et Finance Internationales* de David Guerreiro, Université Paris 8, https://economix.fr/docs/1045/chap_1_2015-16.pdf

20. Site *FOREX* français : <http://www.forex.fr/newslist/8696-la-regle-des-trois-unites-du-marche-des-changes>
21. Investopedia Bid-Ask : <http://www.investopedia.com/terms/b/bid-and-asked.asp>
22. Wikipédia Broker : <https://en.wikipedia.org/wiki/Broker>
23. Wikipédia Régression Statistique : [https://fr.wikipedia.org/wiki/R%C3%A9gression_\(statistiques\)](https://fr.wikipedia.org/wiki/R%C3%A9gression_(statistiques))