

Quelques notes sur l'investissement d'un portefeuille à un actif en présence d'information complémentaire au marché

Thierry Bazier-Matte

10 février 2017

Table des matières

1	Introduction	3
1.1	Avant propos	3
1.2	Exposition du problème et hypothèses	3
1.3	Dimensionnalité de l'information	6
1.4	Risque et garanties statistiques sur la décision	6
1.5	Interprétations	6
1.6	Objectifs	7
2	Optimisation moderne de portefeuille	8
2.1	Théorie classique du portefeuille	8
2.2	Portefeuille universel / Papiers d'Elad Hazan	9
2.3	Théorie de portefeuille régularisé	9
2.4	Fama and French et suivants ?	9
2.5	Articles du NIPS	9
2.6	Papiers de Ben Van Roy	9
2.7	Conclusions : Notre problème par rapport à ces deux disciplines . . .	9
3	Introduction aux fonctions de décisions non linéaires	10
3.1	Introduction	10
3.2	Dualité	11
3.3	Alternate problem	13
4	Garanties statistiques	15
4.1	Bornes de généralisation	15
4.2	Bornes de sous optimalité	17
4.3	Lemmes	18
5	Distributions synthétiques	20

1 Introduction

1.1 Avant propos

[**Todo:** Discuter du rôle croissant que jouent l’informatique et les statistiques dans la construction de portefeuille. Contraster avec les math. stochastiques. Citer Simons et cet article de Quandl selon lequel data is the new shit.]

1.2 Exposition du problème et hypothèses

Ce mémoire vise à établir clairement et rigoureusement comment un investisseur *averse au risque* disposant d’*information complémentaire* au *marché* peut utiliser cette information pour accroître son *utilité espérée* ou, de façon équivalente, son *rendement équivalent certain*.

Modélisation du marché Nous entendrons ici par *marché* n’importe quel type d’actif financier ou spéculatif dans lequel on peut investir une partie de sa fortune dans l’espoir de la voir fructifier au cours d’une période de temps arbitraire. Ainsi, tout au long de l’exposé théorique qui suivra, il peut être pertinent d’avoir en tête les rendements quotidiens issus des grands indices boursiers (par exemple les 500 plus grandes capitalisations américaines). Cependant, le traitement qui sera développé pourrait tout aussi bien s’appliquer à une action cotée en bourse dont on considère les rendements mensuels.^[Nécessaire?] Mathématiquement, l’idée de marché peut ainsi être réduite à celle d’une variable aléatoire $R(t)$ décrivant l’évolution du rendement de l’actif en question.

Relativement à l’idée de marché, nous ferons également l’hypothèse que l’univers a une influence sur ces rendements. Il serait par exemple raisonnable de croire que le prix du pétrole a une influence sur l’évolution du rendement du marché américain. De la même façon, l’annonce d’un scandale aura à son tour des répercussions sur la valeur du titre de la compagnie dont il est l’objet. En outre, il a été montré par Fama et French que le rendement d’une action pouvait s’expliquer comme une combinaison de quelques facteurs fondamentaux (la taille de l’entreprise, le risque de marché et le ratio cours/valeur). On peut alors considérer un vecteur d’information $\vec{X}(t) = (X_1(t), X_2(t), \dots)$ dont chaque composante représente une information particulière, par exemple l’absence ou la présence d’un certain type de scandale, un ratio comptable, le prix d’un certain actif financier.^[Rephrase] D’un point de vue probabiliste, on dira donc qu’il existe une forme de dépendance entre $R(t)$ et $\{\vec{X}(\tau) \mid \tau < t\}$ l’ensemble des événements antérieurs à t . Le processus joint de ces deux événements sera désormais défini comme *la distribution totale de marché*, ou simplement le marché.

Stationarité Bien qu’un tel modèle permette de représenter de façon très générale l’évolution d’un marché, nous formulerons l’hypothèse supplémentaire selon laquelle le marché est un processus *stationnaire*. Ceci permet notamment d’évacuer la notion

temporelle afin de ne représenter qu’une distribution de causes (l’information X) et d’effet (l’observation des rendements R). Cette hypothèse est assez contraignante. Elle suppose d’une part que les réalisations passées n’ont aucun effet sur les réalisations futures (indépendance) et d’autre part que la distribution de marché est figée dans le temps, ce qui implique notamment l’absence de probabilité de faillite. Elle implique aussi que le marché ne peut être vu comme un environnement adversarial qui réagirait par exemple aux décisions d’un investisseur. Ceci vient notamment mettre en cause la théorie des marchés efficients selon laquelle une brèche dans l’absence d’arbitrage serait immédiatement colmatée par des spéculateurs (effet d’autorégulation). Nous aurons toutefois l’occasion de revenir plus en détail sur les liens à faire entre cet exposé et l’efficience des marchés.

Approche mathématique et statistique Dans ce qui suit, nous noterons par M la distribution de marché. Le vecteur aléatoire d’information sera par ailleurs formé de m composantes ; pour l’instant, aucune hypothèse par rapport à la dépendance des composantes de X ne sera formulée. À ce point-ci, on a donc le modèle de marché suivant :

$$M = (R, X_1, \dots, X_m). \quad (1)$$

On fera également l’hypothèse qu’on possède un ensemble de n éléments échantillonnés à partir de M , de sorte que :

$$\{r_i, x_{i1}, \dots, x_{im}\}_{i=1}^n \sim M \quad (2)$$

représente notre ensemble d’échantillonnage (aussi appelé ensemble d’entraînement). Le domaine des rendements possibles de R sera noté $\mathbf{R} \subseteq \mathcal{R}$ et celui du vecteur d’information X sera noté $\mathbf{X} \subseteq \mathcal{R}^m$. Le vecteur d’observations de rendement sera noté $r \in \mathcal{R}^n$ et la matrice d’information par $X \in \mathcal{R}^{n \times m}$.

Modélisation de la préférence Indépendamment de la notion de marché, on a d’autre part l’aspect d’aversion au risque qui est modélisé par une fonction d’utilité $u : \mathbf{R} \rightarrow \mathbf{U}$, où $\mathbf{R} \subseteq \mathcal{R}$ est le domaine (fermé ou non) des rendements considérés et $\mathbf{U} \subseteq \mathcal{R}$ celui des *utilités*.

Bien qu’en pratique il soit plus facile de travailler sur des fonctions possédant des valeurs dans \mathbf{U} , en pratique cet espace est adimensionnel^[Citation needed], de sorte que nos résultats seront présentés dans l’espace des rendements \mathbf{R} .

Fonction de décision Donnés ces éléments de base, le but de ce mémoire sera alors de déterminer une fonction de décision d’investissement $q : \mathbf{X} \rightarrow \mathbf{P} \subseteq \mathcal{R}$ maximisant l’utilité espérée de l’investissement.

Mathématiquement on a donc le problème fondamental suivant :

$$\underset{q \in \mathbf{Q}}{\text{maximiser}} \quad \mathbf{E}u(R \cdot q(X)), \quad (3)$$

où l'optimisation a lieu dans un espace de fonctions \mathcal{Q} à préciser.

Cependant, comme la distribution $(X, R) = M$ est inconnue, il est impossible de déterminer la fonction q^* minimisant cet objectif. On dispose toutefois d'un échantillon de M dont on peut se servir pour approximer le problème (SAA, voir Shapiro^[Citation needed]) :

$$\underset{q \in \mathcal{Q}}{\text{maximiser}} \quad n^{-1} \sum_{i=1}^n u(r_i q(x_i)), \quad (4)$$

mais encore ici le problème est mal spécifié, puisqu'aucune contrainte n'a été posée sur l'espace \mathcal{Q} . Par exemple, il suffirait de prendre pour q un dictionnaire associant à x_i la valeur αr_i , où $\alpha > 0$, et à toute autre valeur de x une valeur nulle pour avoir une valeur d'utilité arbitrairement grande à mesure que $\alpha \rightarrow \infty$.

Risque in-échantillon et hors échantillon une telle fonction q est qu'elle se généralise très mal. En effet pour toute observation x qui ne figurerait dans l'ensemble d'entraînement, q prescrirait alors un investissement nul. Il y a alors une énorme différence entre l'utilité observée au sein de notre échantillon et l'utilité hors échantillon.

Donnée une fonction de décision $q \in \mathcal{Q}$ et un échantillon de M , on définit le *risque in-échantillon* ou *risque empirique* par

$$\hat{R}(q) = n^{-1} \sum_{i=1}^n \ell(r_i q(x_i)), \quad (5)$$

où $\ell = -u$. De la même façon, on définit le *risque hors-échantillon* ou *erreur de généralisation* par

$$R(q) = \mathbf{E} \ell(R \cdot q(X)). \quad (6)$$

On peut souhaiter d'une bonne fonction de décision qu'elle performe bien hors échantillon, aussi la quantité $R(q) - \hat{R}(q)$ sera-t-elle primordiale et beaucoup d'attention lui sera consacrée dans les prochaines sections. Notons que le risque hors-échantillon étant théoriquement impossible à calculer, en pratique on segmentera l'ensemble d'échantillonnage en deux parties, l'une dédiée à l'apprentissage, l'autre à évaluer la performance hors échantillon.

Régularisation Afin de contrecarrer le risque hors échantillon, la solution est en fait de pénaliser la complexité de la fonction de décision q (rasoir d'Occam). Ainsi, on étudiera en profondeur le choix d'une fonctionnelle $R : \mathcal{Q} \rightarrow \mathcal{R}$ permettant de quantifier la complexité de q . L'objectif serait alors

$$\underset{q \in \mathcal{Q}}{\text{maximiser}} \quad n^{-1} \sum_{i=1}^n u(r_i q(x_i)) - R(q). \quad (7)$$

Par exemple, comme les mesures sur x peuvent comporter de l'incertitude ou du bruit, il serait souhaitable que la décision $q(x_1)$ soit proche de $q(x_2)$, si x_1 et x_2 sont eux

même proches dans l'espace \mathbf{X} . Si R encodait une telle préférence, ne fonction discontinue comme le dictionnaire présenté plus haut sera alors hautement défavorisée, et une fonction plus lisse y serait préférée.

[**Todo:** Introduire la validation croisée ainsi que le paramètre λ dans l'objectif.]

Espaces de décision En pratique, ce mémoire ne considérera que des espaces de Hilbert pour \mathbf{Q} . Un des avantages des espaces de Hilbert, c'est qu'ils induisent naturellement une notion de norme $\|\cdot\|_H$, qu'on peut intuitivement relier au concept de complexité. Nous nous intéresserons donc aux propriétés induites par $R(q) = \|q\|_H^2 = \langle q, q \rangle$. Il y a aussi moyen, sous des conditions assez techniques (théorème de la représentation) de généraliser la norme L_2 de q à une norme L_p général. En particulier, nous verrons qu'une régularisation donnée par norme L_1 induit certaines propriétés d'éparsité dans la solution.

Décisions linéaires De façon générale, la forme de décision la plus simple est celle qui combine linéairement les p observations de $x \in \mathbf{X} \subseteq \mathcal{R}^p$; autrement dit lorsque qu'on contraint $\mathbf{Q} = \mathbf{X}^*$, i.e., à l'espace dual de \mathbf{X} . En langage plus clair, à toute fonction $q \in \mathbf{Q}$ il existe un vecteur de dimension p tel que la décision dérivée de l'observation x sera donnée par $q(x) = \langle q, x \rangle = q^T x$.

La régularisation L_2 de q devient alors tout simplement $R(q) = q^T q = \|q\|^2$ et la fonction optimale de décision q^* sera alors déterminée en résolvant le problème d'optimisation suivant :

$$\underset{q \in \mathcal{R}^p}{\text{maximiser}} \quad n^{-1} \sum_{i=1}^n u(r_i q^T x_i) - \lambda \|q\|^2. \quad (8)$$

1.3 Dimensionnalité de l'information

[**Todo:** Discussion du phénomène big data, de l'importance de p]

1.4 Risque et garanties statistiques sur la décision

[**Todo:** Discussion sur les méthodes de risques hors échantillon, complexité de l'échantillonnage, mesure Rademacher, distance par rapport à la "meilleure" décision]

1.5 Interprétations

Interprétation géométrique dans l'espace \mathbf{X}

Interprétation statistique (avec matrix covariance)

Autre ?

1.6 Objectifs

2 Optimisation moderne de portefeuille

Dans ce document, nous allons tenter de classer et de répertorier la plupart des méthodes ayant rapport, de près ou de loin, à l'intersection des méthodes statistiques avancées et de l'apprentissage machine avec la théorie du portefeuille, en présentant pour chacune d'elle leurs avantages et leurs inconvénients.

2.1 Théorie classique du portefeuille

Une revue de littérature sur la théorie du portefeuille serait fondamentalement incomplète sans l'article fondateur de Markowitz, publié en 1952 [?].

Nous allons montrer que le cadre théorique développé par Markowitz peut être considéré comme un cas particulier de notre algorithme, pour autant que l'on considère un portefeuille à un seul actif.

Soit $w \in \mathcal{R}^k$ le vecteur représentant la répartition du portefeuille de Markowitz à k actifs à optimiser. Alors un investisseur *markowitzien* souhaite résoudre le problème suivant :

$$\begin{aligned} &\text{minimiser} && w^T \Sigma w \\ &\text{tel que} && \mu^T w = \mu_0, \end{aligned} \tag{9}$$

où $\Sigma \in \mathcal{R}^{k \times k}$ est la covariance du rendement des actifs et $\mu \in \mathcal{R}^k$ le vecteur d'espérance. **[Todo: Montrer formellement.]** Par la théorie de l'optimisation convexe, il existe une constante $\gamma \in \mathcal{R}$ telle que le problème énoncé est équivalent à

$$\text{maximiser} \quad \mu^T w + \gamma w^T \Sigma w. \tag{10}$$

Dans le cas où on considère un portefeuille à un seul actif, alors ce problème se réduit alors à

$$\text{maximiser} \quad \mu q - \gamma \sigma^2 q^2, \tag{11}$$

où on a posé $\mu := \mathbf{E}R$ et $\sigma^2 := \text{Var } R$.

Supposons qu'un investisseur soit doté d'une utilité quadratique paramétrée par

$$u(r) = r - \frac{\gamma}{\sigma^2 + \mu^2} \sigma^2 r^2, \tag{12}$$

et que l'information factorielle intégrée à l'algorithme ne consiste uniquement qu'en les rendements eux mêmes ; autrement dit, le vecteur d'information X se réduirait tout simplement à un terme constant fixé à 1, *i.e.*, $X \sim 1$. **[Todo: expliquer].**

Avec une utilité (12) et l'absence d'information supplémentaire, l'objectif de (??) devient aussitôt

$$\mathbf{E}U(qR) = q\mathbf{E}R - \frac{\gamma}{\sigma^2 + \mu^2} \sigma^2 q^2 \mathbf{E}R^2. \tag{13}$$

Or, puisque $\text{Var } R = \mathbf{E}R^2 - (\mathbf{E}R)^2$, on a que $\mathbf{E}R^2 = \sigma^2 + \mu^2$, ce qui entraîne donc que (12) s'exprime par

$$\text{maximiser } \mathbf{E}U(qR) = \mu q - \gamma \sigma^2 q^2, \quad (14)$$

ce qui est tout à fait identique à (11).

Nous suggérons au lecteur intéressé par l'équivalence des diverses formulations d'optimisation de portefeuille dans un univers de Markowitz [?] et [?], tous deux publiés à l'occasion du soixantième anniversaire de [?].

2.2 Portefeuille universel / Papiers d'Elad Hazan

Ce mémoire sera également consacré aux garanties statistiques de performance des estimateurs q^* .

Bien que le modèle soit différent et de nature itérative, le *portefeuille universel* de [?] est à notre connaissance un des premiers modèles de gestion de portefeuille à exploiter une distribution arbitraire tout en proposant des garanties statistiques de convergence.

Voir [?, ?].

2.3 Théorie de portefeuille régularisé

[?]

2.4 Fama and French et suivants ?

[?]

2.5 Articles du NIPS

2.6 Papiers de Ben Van Roy

2.7 Conclusions : Notre problème par rapport à ces deux disciplines

3 Introduction aux fonctions de décisions non linéaires

3.1 Introduction

Soit $\kappa : \mathbf{X} \times \mathbf{X} \rightarrow \mathcal{R}$ un noyau semi-défini positif, \mathbf{H} l'espace de Hilbert à noyau reproduisant induit par κ et $K \in \mathcal{R}^{n \times n}$ la matrice associée à l'ensemble d'échantillonnage $S_n \sim M^n$. Le problème d'optimisation de portefeuille régularisé s'exprime alors par

$$\underset{q \in \mathbf{H}}{\text{maximiser}} \quad n^{-1} \sum_{i=1}^n u(r_i q(x_i)) - \lambda \|q\|_{\mathbf{H}}^2. \quad (15)$$

Tel que mentionné, la dimension de \mathbf{H} est possiblement infinie, ce qui rend numériquement impossible la recherche d'une solution q^* . Toutefois, le théorème de la représentation permet de rendre le problème résoluble.

Théorème. *Toute solution q^* de (15) repose dans le sous-espace vectoriel engendré par l'ensemble des n fonctions $\{\phi_i\}$, où $\phi_i = \kappa(x_i, \cdot)$. Numériquement, il existe un vecteur $\alpha^* \in \mathcal{R}^n$ tel que,*

$$q^* = \sum_{i=1}^n \alpha_i^* \phi_i = (\alpha^*)^T \phi. \quad (16)$$

Démonstration. Voir [?], Théorème 5.4 pour une démonstration tenant compte d'un objectif régularisé général. La démonstration est due à [?]. \square

Le théorème de la représentation permet donc de chercher une solution dans un espace à n dimensions, plutôt que la dimension possiblement infinie de \mathbf{H} . En effet, puisque

$$q^* = \sum_{i=1}^n \alpha_i^* \phi_i, \quad (17)$$

où $\alpha \in \mathcal{R}^n$ [**Todo:** Espace cotangent ???], on peut donc restreindre le domaine d'optimisation à \mathcal{R}^n . L'objectif de (15) devient alors

$$n^{-1} \sum_{i=1}^n u(r_i \sum_{j=1}^n \alpha_j \phi_j(x_i)) - \lambda \langle q, q \rangle_{\mathbf{H}}. \quad (18)$$

Le premier terme se réexprime comme

$$n^{-1} \sum_{i=1}^n u(r_i \alpha^T \phi(x_i)), \quad (19)$$

alors qu'en employant les propriétés de linéarité du produit intérieur, on transforme le second terme par

$$\langle q, q \rangle_{\mathbf{H}} = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle \phi_i, \phi_j \rangle_{\mathbf{H}} \quad (20)$$

$$= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \kappa(x_i, x_j) \quad (21)$$

$$= \alpha^T K \alpha. \quad (22)$$

De sorte que le problème général (15) peut se reformuler par

$$\boxed{\text{maximiser}_{\alpha \in \mathcal{R}^n} \quad n^{-1} \sum_{i=1}^n u(r_i \alpha^T \phi(x_i)) - \lambda \alpha^T K \alpha} \quad (23)$$

3.2 Dualité

Let us consider the following problem, optimized over $q \in \mathcal{R}^p$:

$$\text{minimiser} \quad \sum_{i=1}^n \ell(r_i q^T x_i) + n\lambda \|q\|^2, \quad (24)$$

where $\ell = -u$. Alternatively, this problem can be respecified using slack vector $\xi \in \mathcal{R}^n$ as

$$\begin{aligned} &\text{minimiser} \quad \sum_{i=1}^n \ell(\xi_i) + n\lambda \|q\|^2 \\ &\text{tel que} \quad \xi_i = r_i q^T x_i. \end{aligned} \quad (25)$$

Let $\alpha \in \mathcal{R}^n$. The Lagrangian of (25) can be written as

$$\mathcal{L}(q, \xi, \alpha) = \sum_{i=1}^n \ell(\xi_i) + n\lambda \|q\|^2 + \sum_{i=1}^n \alpha_i (r_i q^T x_i - \xi_i). \quad (26)$$

Because the objective (25) is convex and its constraints are affine in q and ξ , Slater's theorem states that the duality gap of the problem is zero. In other words, solving (24) is equivalent to maximizing the Lagrange dual function g over α :

$$\text{maximiser} \quad g(\alpha) = \inf_{q, \xi} \mathcal{L}(q, \xi, \alpha). \quad (27)$$

Now, note that

$$g(\alpha) = \inf_{q, \xi} \left\{ \sum_{i=1}^n \ell(\xi_i) + n\lambda \|q\|^2 + \sum_{i=1}^n \alpha_i (r_i q^T x_i - \xi_i) \right\} \quad (28)$$

$$= \inf_{\xi} \left\{ \sum_{i=1}^n \ell(\xi_i) - \alpha^T \xi \right\} + \inf_q \left\{ \sum_{i=1}^n \alpha_i r_i q^T x_i + n\lambda \|q\|^2 \right\} \quad (29)$$

$$= -\sup_{\xi} \left\{ \alpha^T \xi - \sum_{i=1}^n \ell(\xi_i) \right\} + \inf_q \left\{ \sum_{i=1}^n \alpha_i r_i q^T x_i + n\lambda \|q\|^2 \right\} \quad (30)$$

$$= - \sum_{i=1}^n \ell^*(\alpha_i) + \inf_q \left\{ \sum_{i=1}^n \alpha_i r_i q^T x_i + n\lambda \|q\|^2 \right\}. \quad (31)$$

Where ℓ^* is the convex conjugate of the loss function and is defined by

$$\ell(\alpha_i) = \sup_{\xi_i} \{ \alpha_i \xi_i - \ell(\xi_i) \}. \quad (32)$$

Note that the identity

$$f(\xi_1, \dots, \xi_n) = \sum_{i=1}^n \ell(\xi_i) \implies f^*(\xi_1, \dots, \xi_n) = \sum_{i=1}^n \ell^*(\xi_i) \quad (33)$$

was used. Consider now the second part of (31). Since the expression is differentiable, we can analytically solve for q :

$$\nabla_q \left\{ \sum_{i=1}^n \alpha_i r_i q^T x_i + n\lambda \|q\|^2 \right\} = 0 \quad (34)$$

implies that

$$q = -\frac{1}{2n\lambda} \sum_{i=1}^n \alpha_i r_i x_i \quad (35)$$

at the infimum.

Using (35), we can eliminate q from (31), so that

$$g(\alpha) = - \sum_{i=1}^n \ell^*(\alpha_i) - \frac{1}{2n\lambda} \sum_{i,j=1}^n \alpha_i \alpha_j r_i r_j x_i^T x_j + \frac{1}{4n\lambda} \sum_{i,j=1}^n \alpha_i \alpha_j r_i r_j x_i^T x_j \quad (36)$$

$$= - \sum_{i=1}^n \ell^*(\alpha_i) - \frac{1}{4n\lambda} (\alpha \circ r)^T K (\alpha \circ r). \quad (37)$$

Therefore, in its dual form, the problem (24) is equivalent to solving

$$\text{minimiser } \sum_{i=1}^n \ell^*(\alpha_i) + \frac{1}{4n\lambda} (\alpha \circ r)^T K (\alpha \circ r). \quad (38)$$

Prescribed investment In its original form, given a feature vector \tilde{x} , the algorithm (24) suggests an investment size of $p_0 = q^T \tilde{x}$, where q is the trained value obtained by optimizing (24). In the dual formulation (38), with optimal value α , we have from (35) :

$$p_0 = q^T x_0 \quad (39)$$

$$= -\frac{1}{2n\lambda} \sum_{i=1}^n \alpha_i r_i x_i^T x_0. \quad (40)$$

[**Todo:** Insert kernel formulation with vector ϕ .]

3.3 Alternate problem

We now consider a new problem, slightly different from (24) where a regularization based on the sum of the square of the investment sizes $q^T x_i$ is applied :

$$\text{minimiser } \sum_{i=1}^n \ell(r_i q^T x_i) + \gamma \sum_{i=1}^n (q^T x_i)^2 + n\lambda \|q\|^2. \quad (41)$$

Again, this problem can be respecified using slack vector $\xi \in \mathcal{R}^n$ as

$$\begin{aligned} \text{minimiser } & \sum_{i=1}^n \ell(\xi_i) + \gamma \sum_{i=1}^n (\xi_i/r_i)^2 + n\lambda \|q\|^2 \\ \text{tel que } & \xi_i = r_i q^T x_i. \end{aligned} \quad (42)$$

The constraints in (42) are again affine, so that Slater's theorem apply.

The lagrangian of (42) is

$$\mathcal{L}(q, \xi, \alpha) = \sum_{i=1}^n \ell(\xi_i) + \gamma \sum_{i=1}^n (\xi_i/r_i)^2 + n\lambda \|q\|^2 + \sum_{i=1}^n \alpha_i (r_i q^T x_i - \xi_i), \quad (43)$$

and we seek its infimum over (q, ξ) .

$$\inf_{q, \xi} \left\{ \sum_{i=1}^n \ell(\xi_i) + \gamma \sum_{i=1}^n (\xi_i/r_i)^2 + n\lambda \|q\|^2 + \sum_{i=1}^n \alpha_i (r_i q^T x_i - \xi_i) \right\} \quad (44)$$

$$= \inf_{\xi} \left\{ \sum_{i=1}^n \ell(\xi_i) + \gamma \sum_{i=1}^n (\xi_i/r_i)^2 - \alpha^T \xi \right\} + \inf_q \left\{ \sum_{i=1}^n \alpha_i r_i q^T x_i - n\lambda \|q\|^2 \right\} \quad (45)$$

$$= -\sup_{\xi} \left\{ \alpha^T \xi - \left(\sum_{i=1}^n \ell(\xi_i) + \gamma \sum_{i=1}^n (\xi_i/r_i)^2 \right) \right\} - \frac{1}{4n\lambda} (\alpha \circ r)^T K (\alpha \circ r). \quad (46)$$

Let $f_i(\xi_i) := h_1(\xi_i) + h_2(\xi_i) = \ell(\xi_i) + \gamma(\xi_i/r_i)^2$. Then, using (33), the first expression of (46) can be restated as

$$-\sup_{\xi} \left\{ \alpha^T \xi - \sum_{i=1}^n f_i(\xi_i) \right\} = -\sum_{i=1}^n f_i^*(\alpha_i). \quad (47)$$

Let us introduce another identity :

$$(h_1 + h_2)^*(\alpha_i) = \inf_{\alpha'_i + \alpha''_i = \alpha_i} \{h_1^*(\alpha'_i) + h_2^*(\alpha''_i)\}. \quad (48)$$

Using (48), (47) can be written as

$$-\sum_{i=1}^n f_i^*(\alpha_i) = -\sum_{i=1}^n (h_1 + h_2)^*(\xi_i) \quad (49)$$

$$= - \sum_{i=1}^n \inf_{\alpha'_i + \alpha''_i = \alpha_i} \{h_1^*(\alpha'_i) + h_2^*(\alpha''_i)\}. \quad (50)$$

The first conjugate function h_1^* is simply ℓ^* . The second conjugate function can be derived analytically :

$$h_2^*(\alpha''_i) = \sup_{\xi_i} \{\alpha''_i \xi_i - h_2(\xi_i)\} \quad (51)$$

$$= \sup_{\xi_i} \{\alpha''_i \xi_i - \gamma(\xi_i/r_i)^2\}. \quad (52)$$

The supremum occurs when

$$\xi_i = \frac{r_i^2}{2\gamma} \alpha''_i. \quad (53)$$

Therefore, (52) simplifies to

$$h_2^*(\alpha''_i) = \frac{r_i^2}{4\gamma} (\alpha''_i)^2. \quad (54)$$

Putting it all back together, the dual of (41) is

$$- \sum_{i=1}^n \inf_{\alpha'_i + \alpha''_i = \alpha_i} \left\{ \ell^*(\alpha'_i) + \frac{r_i^2}{4\gamma} (\alpha''_i)^2 \right\} - \frac{1}{4n\lambda} (\alpha \circ r)^T K(\alpha \circ r), \quad (55)$$

which is equivalent to

$$- \sum_{i=1}^n \ell^*(\alpha_i) - \frac{1}{4\gamma} \sum_{i=1}^n (r_i \beta_i)^2 - \frac{1}{4n\lambda} (r \circ (\alpha + \beta))^T K(r \circ (\alpha + \beta)), \quad (56)$$

with new optimization variables $\alpha = \alpha', \beta = \alpha'' \in \mathcal{R}^n$. The dual optimization problem is therefore

$$\text{minimiser} \quad \sum_{i=1}^n \ell^*(\alpha_i) + \frac{1}{4\gamma} \|r \circ \beta\|^2 + \frac{1}{4n\lambda} (r \circ (\alpha + \beta))^T K(r \circ (\alpha + \beta)). \quad (57)$$

Prescribed investment [Todo:]

4 Garanties statistiques

4.1 Bornes de généralisation

Exposition du problème Soit \mathcal{Q} un espace de Hilbert à noyau reproduisant induit par κ et soit un ensemble d'entraînement $\mathcal{S}_n = \{(x_i, y_i)\}_{i=1}^n \sim M^n$ échantillonné à partir de la distribution de marché. Alors on peut définir l'*algorithme de décision* $\mathcal{Q} : M^n \rightarrow \mathcal{Q}$ par

$$\mathcal{Q}(\mathcal{S}_n) = \arg \max_{q \in \mathcal{Q}} \left\{ \widehat{EU}(\mathcal{S}_n, q) - \lambda \|q\|^2 \right\}. \quad (58)$$

Comme on l'a vu, résoudre (58) est aussi équivalent à

$$\underset{\alpha \in \mathcal{R}^n}{\text{maximiser}} \quad n^{-1} \sum_{i=1}^n u(r_i \alpha^T \phi(x_i)) - \lambda \|\alpha\|_K^2, \quad (59)$$

où $\phi : \mathcal{R}^p \rightarrow \mathcal{R}^n$ le vecteur d'application induit par la matrice d'information Ξ . La relation $q = \alpha^T \phi$ permet de passer d'une représentation à l'autre.

La question qui se pose naturellement est de savoir dans quelle mesure une fonction de décision $\hat{q} = \mathcal{Q}(\mathcal{S}_n)$ est capable d'offrir à un investisseur une utilité espérée comparable à celle qu'il aurait observée au sein de l'ensemble d'entraînement. Il serait aussi souhaitable qu'une telle garantie soit indépendante de l'ensemble d'entraînement \mathcal{S}_n . Autrement dit, on cherche à déterminer une borne probabiliste Ω sur l'erreur de généralisation de $\hat{q} = \mathcal{Q}(\mathcal{S}_n)$ valide pour tout $\mathcal{S}_n \sim M^n$:

$$\hat{\zeta}(\mathcal{S}_n) \leq \Omega(n, \dots), \quad (60)$$

où

$$\hat{\zeta}(\mathcal{S}_n) = \widehat{EU}(\mathcal{S}_n, \mathcal{Q}(\mathcal{S}_n)) - EU(\mathcal{Q}(\mathcal{S}_n)) \quad (61)$$

représente l'erreur de généralisation. On peut en fait montrer que $\Omega \rightarrow 0$ à mesure que $n \rightarrow \infty$. En fait, on peut démontrer que $\Omega = O(n^{-1/2})$, ce qui permet de quantifier la "vitesse" à laquelle la convergence à lieu.

Démonstration Considérons deux ensembles d'entraînement : $\mathcal{S}_n \sim M^n$ et \mathcal{S}'_n , où \mathcal{S}'_n ne diffère de \mathcal{S}_n que par un seul point (par exemple le j -ème point serait rééchantillonné de la distribution de marché M). De l'algorithme \mathcal{Q} on dérivera alors deux décisions : \hat{q} et \hat{q}' . Pour n suffisamment grand, on peut alors s'attendre à ce que l'utilité dérivée de ces deux décisions soit relativement proche, et ce, pour toute observation. On aurait alors une borne $\beta(n)$ telle que pour tout $(x, r) \sim M$,

$$|u(r \hat{q}(x)) - u(r \hat{q}'(x))| \leq \beta. \quad (62)$$

C'est ce qu'on appelle dans la littérature la *stabilité algorithmique*. La plupart des algorithmes régularisés classiques disposent par ailleurs d'une telle stabilité. En particulier,

le terme de régularisation $\lambda\|q\|^2$, combiné à la continuité Lipschitz de u font en sorte que $\beta = O(n^{-1})$.

Doté de cette stabilité de \mathcal{Q} , on peut alors borner la différence dans l'erreur de généralisation de \mathcal{S}_n et \mathcal{S}'_n :

$$|\hat{\zeta}(\mathcal{S}_n) - \hat{\zeta}(\mathcal{S}'_n)| = |\mathbf{EU}(\hat{q}) - \mathbf{EU}(\hat{q}') + \widehat{\mathbf{EU}}(\mathcal{S}_n, \hat{q}) - \widehat{\mathbf{EU}}(\mathcal{S}'_n, \hat{q}')| \quad (63)$$

$$\leq |\mathbf{EU}(\hat{q}) - \mathbf{EU}(\hat{q}')| + |\widehat{\mathbf{EU}}(\mathcal{S}_n, \hat{q}) - \widehat{\mathbf{EU}}(\mathcal{S}'_n, \hat{q}')|. \quad (64)$$

Or, par le théorème de Jensen appliqué à la fonction valeur absolue, on obtient du premier terme que

$$|\mathbf{EU}(\hat{q}) - \mathbf{EU}(\hat{q}')| = |\mathbf{E}(u(R \cdot \hat{q}(X)) - u(R \cdot \hat{q}'(X)))| \quad (65)$$

$$\leq \mathbf{E}(|u(R \cdot \hat{q}(X)) - u(R \cdot \hat{q}'(X))|) \quad (66)$$

$$\leq \beta, \quad (67)$$

par définition de la stabilité. Quant au deuxième terme de (64) on peut le borner de la même façon :

$$|\widehat{\mathbf{EU}}(\mathcal{S}_n, \hat{q}) - \widehat{\mathbf{EU}}(\mathcal{S}'_n, \hat{q}')| \quad (68)$$

$$= n^{-1} \left| \sum_{i=1}^n \mathbb{I}_{i \neq j} u(r_i \hat{q}(x_i)) + u(r_j \hat{q}(x_j)) - \sum_{i=1}^n \mathbb{I}_{i \neq j} u(r_i \hat{q}'(x_i)) - u(r'_j \hat{q}'(x'_j)) \right| \quad (69)$$

$$\leq n^{-1} \left(|u(r_j \hat{q}(x_j)) - u(r'_j \hat{q}'(x'_j))| + \sum_{i=1}^n \mathbb{I}_{i \neq j} |u(r_i \hat{q}(x_i)) - u(r_i \hat{q}'(x_i))| \right) \quad (70)$$

$$\leq n^{-1} (|u(r_j \hat{q}(x_j)) - u(r'_j \hat{q}'(x'_j))| + (n-1)\beta). \quad (71)$$

Considérons le premier terme. Par le lemme **[Todo:]**, On sait que $\hat{q}(x) \leq (2\lambda)^{-1} \bar{r} \xi^2$ et que $|R| \leq \bar{r}$. On peut donc borner cette différence par la différence dans l'utilité dérivée par la meilleure décision d'investissement sur le meilleur rendement et sur le pire rendement. Par hypothèse Lipschitz et de sous-gradient de 1 à $r = 0$, on sait que pour $r > 0$, $u(r) < r$ et que pour $r < 0$, $\gamma r \leq u(r)$. On peut donc conclure que

$$|u(r_j \hat{q}(x_j)) - u(r'_j \hat{q}'(x'_j))| \leq u((2\lambda)^{-1} \bar{r}^2 \xi^2) - u(-(2\lambda)^{-1} \bar{r}^2 \xi^2) \quad (72)$$

$$\leq (2\lambda)^{-1} (\gamma + 1) \bar{r}^2 \xi^2. \quad (73)$$

Ce qui entraîne donc que

$$|\widehat{\mathbf{EU}}(\mathcal{S}_n, \hat{q}) - \widehat{\mathbf{EU}}(\mathcal{S}'_n, \hat{q}')| \leq \frac{\gamma + 1}{2\lambda n} \bar{r}^2 \xi^2 + \frac{n-1}{n} \beta \quad (74)$$

$$\leq \beta + \frac{\gamma + 1}{2\lambda n} \bar{r}^2 \xi^2. \quad (75)$$

Next McDiarmid et autre. Easy stuff.

Équivalent certain Puis inverser pour obtenir l'équivalent certain.

Note bibliographique La théorie de la stabilité algorithmique remonte en fait aux années 70 avec les travaux de Luc Devroye appliqués à l'algorithme des k plus proches voisins^[Citation needed]. Jusqu'alors, les bornes de généralisation étaient présentées pour toute décision $q \in \mathcal{Q}$ (ie Vapnik). Bousquet^[Citation needed] a été le premier à présenter des résultats dans des espaces de Hilbert à noyau reproduisant. La démonstration est fortement inspirée de l'excellente référence Mohri^[Citation needed]. La démonstration de la borne de la décision bornée est un résultat inédit, dû à Delage dans le cas linéaire.

4.2 Bornes de sous optimalité

Exposition du problème Jusqu'ici, les efforts théoriques ont été déployés pour déterminer comment se comportait la fonction de décision $\hat{q} = \mathcal{Q}(\mathcal{S}_n)$ dans un univers probabiliste par rapport à l'univers statistique dans lequel elle avait été construite. Notre attention va maintenant se tourner vers la performance de \hat{q} dans l'univers probabiliste par rapport à la meilleure décision disponible, c'est à dire la solution q^* de

$$\underset{q \in \mathcal{Q}}{\text{maximiser}} \quad \mathbf{E}U(R \cdot q(X)). \quad (76)$$

Il convient cependant de réaliser que l'existence d'une borne sur q^* n'est pas assurée. En effet, supposons d'une part que l'on dispose d'une utilité neutre au risque, telle que $u(r) = r$, et d'autre part que $\mathbf{E}R = 0$. Soit $\alpha > 0$. On pourrait alors définir la fonction suivante :

$$q = \alpha \mathbf{E}(R \kappa(X, \cdot)) \quad (77)$$

On aurait alors

$$\mathbf{E}U(q) = \mathbf{E}(Rq(X)) = \mathbf{E}(R\mathbf{E}(R \kappa(X, X))) \quad (78)$$

$$= \mathbf{E}(R^2 \kappa(X, X)) \geq 0, \quad (79)$$

On peut alors obtenir une utilité espérée non bornée à mesure que $\alpha \rightarrow \infty$. Par ailleurs, ainsi défini, q représente effectivement la covariance entre R et la projection de X dans l'espace dual de \mathcal{Q} . Puisque l'utilité est neutre, on sait qu'en espérance l'application de q à X variera de la même façon que celle de R et donc qu'on aura une utilité infinie. On verra plus loin au cours d'une démonstration la motivation derrière cette hypothèse supplémentaire :

Hypothèse 1. *L'utilité croît sous-linéairement, ie. $u(r) = o(r)$.*

Une autre hypothèse est maintenant nécessaire pour s'assurer que q^* soit borné : l'efficience des marchés. Dans notre cadre théorique, ceci se traduit par l'absence de l'existence d'une fonction $q \in \mathcal{Q}$ telle que

$$\mathbf{P}\{R \cdot q(X) > 0\} = 1. \quad (80)$$

D'un point de vue strictement financier, cela fait certainement du sens en vertu de l'efficience des marchés, version semi-forte^[Citation needed]. D'un point de vue théorique, ceci exige en fait qu'il n'y ait pas de région dans \mathbf{X} telle que tous les rendements s'y produisant soient nécessairement positifs ou négatifs.**[Todo: Insérer image].**

Hypothèse 2. Pour toute région $\mathcal{R} \subseteq \mathbf{X}$,

$$P\{R \geq 0 \mid X \in \mathcal{R}\} < 1, \quad (81)$$

et de la même façon avec l'évènement $P\{R \leq 0\}$.

Borne On cherchera donc à établir une borne sur l'erreur de sous-optimalité de $\hat{q} \sim \mathcal{Q}(M^n)$.

4.3 Lemmes

Stabilité On montre ici que

$$\beta \leq \frac{(\gamma \bar{r} \xi)^2}{2\lambda n}. \quad (82)$$

Borne sur la décision algorithmique On va ici démontrer que la décision $\hat{q}(x)$ est bornée, et ce, pour tout $x \in \mathbf{X}$ et pour toute solution \hat{q} de

$$\underset{q \in \mathcal{Q}}{\text{maximiser}} \quad n^{-1} \sum_{i=1}^n u(r_i q(x_i)) - \lambda \|q\|^2. \quad (83)$$

Pour ce faire, on va mettre à profit la propriété reproductive de \mathcal{Q} induite par κ . En effet, celle-ci stipule que

$$q(x) = \langle q, \kappa(x, \cdot) \rangle_{\mathcal{Q}} \leq \|q\| \sqrt{\kappa(x, x)}, \quad (84)$$

où l'inégalité découle de l'inégalité Cauchy-Schwartz appliquée au produit interne de \mathcal{Q} . On rappelle que, par hypothèse, $\forall x \in \mathbf{X}, \kappa(x, x) \leq \xi^2$; il suffit donc de borner $\|q\|_{\mathcal{Q}}$. Or, puisque $u(r) \leq r$, on remarque que

$$n^{-1} \sum_{i=1}^n u(r_i q(x_i)) \leq n^{-1} \sum_{i=1}^n r_i q(x_i) \quad (85)$$

$$\leq n^{-1} \sum_{i=1}^n r_i \sqrt{\kappa(x_i, x_i)} \|q\| \quad (86)$$

$$\leq \bar{r} \xi \|q\|. \quad (87)$$

Puisque l'expression $\bar{r} \xi \|q\| - \lambda \|q\|^2$ est quadratique et atteint son maximum lorsque

$$\|q\| = \frac{\bar{r} \xi}{2\lambda}, \quad (88)$$

on en conclut que $\|\hat{q}\| \leq (2\lambda)^{-1}\bar{r}\xi$ et donc que

$$\hat{q}(x) \leq \frac{\bar{r}\xi^2}{2\lambda}. \quad (89)$$

[**Todo:** Montrer qu'on peut effectivement montrer que la borne de l'expression est plus grande que celle du problème initial... Pour ce faire, utiliser a) le sous gradient de u et b) la dominance de q' sur q .]

Forte concavité L'objectif est fortement concave, que ce soit sous sa version statistique \widehat{EU}_λ ou probabiliste EU_λ . Autrement dit, pour tout $\alpha \in [0, 1]$, on a

$$EU_\lambda(q_1 + (1-\alpha)q_2) \geq \alpha EU_\lambda(q_1) + (1-\alpha)EU_\lambda(q_2) + \lambda\alpha(1-\alpha)\|q_1 - q_2\|^2, \quad (90)$$

et de même pour \widehat{EU}_λ . Effectivement, puisque u est concave et $\|\cdot\|^2$ est convexe, on a successivement :

$$EU_\lambda(\alpha q_1 + (1-\alpha)q_2) \quad (91)$$

$$= Eu(R \cdot (\alpha q_1 + (1-\alpha)q_2)(X)) - \lambda\|\alpha q_1 + (1-\alpha)q_2\|^2 \quad (92)$$

$$= Eu(\alpha(R \cdot q_1(X)) + (1-\alpha)(R \cdot q_2(X))) - \lambda\|\alpha q_1 + (1-\alpha)q_2\|^2 \quad (93)$$

$$\geq E(\alpha u(R \cdot q_1(X)) + (1-\alpha)u(R \cdot q_2(X))) - \lambda\|\alpha q_1 + (1-\alpha)q_2\|^2 \quad (94)$$

$$= \alpha EU(q_1) + (1-\alpha)EU(q_2) - \lambda\|\alpha q_1 + (1-\alpha)q_2\|^2 \quad (95)$$

$$= \alpha EU_\lambda(q_1) + (1-\alpha)EU_\lambda(q_2) - \lambda(\|\alpha q_1 + (1-\alpha)q_2\|^2 - \alpha\|q_1\|^2 - (1-\alpha)\|q_2\|^2) \quad (96)$$

$$\geq \alpha EU_\lambda(q_1) + (1-\alpha)EU_\lambda(q_2) - \lambda(\alpha\|q_1\|^2 + (1-\alpha)\|q_2\|^2 - \alpha\|q_1\|^2 - (1-\alpha)\|q_2\|^2) \quad (97)$$

$$= \alpha EU_\lambda(q_1) + (1-\alpha)EU_\lambda(q_2). \quad (98)$$

La preuve est la même lorsqu'on considère \widehat{EU}_λ .

Borne sur la décision optimale On veut montrer que $\|q^*\|$ est borné. Pour ce faire, on va tout d'abord décomposer $q = s\theta$, où on pose $\|\theta\| = 1$ et $s > 0$; ainsi on peut poser notre problème d'optimisation comme la recherche d'une 'direction' θ et d'une magnitude s dans \mathcal{Q} . De plus, puisque $\|q\| = s$, il suffit de montrer que s^* est borné.

Notons d'abord que l'hypothèse 2 entraîne en particulier qu'il existe $\delta > 0$ tel que

$$P\{R \cdot \theta(X) \leq -\delta\} > \varrho \geq 0 \quad (99)$$

pour tout $\theta \in \mathcal{Q}$ tel que $\|\theta\| = 1$. Définissons maintenant une variable aléatoire à deux états : $B = -\delta$ avec probabilité ϱ et $B = \bar{r}\xi$ avec probabilité $1 - \varrho$. Puisque $R \cdot \theta(X) \leq \bar{r}\xi$, on a alors que, pour tout $r \in \mathbf{R}$,

$$P\{B \geq r\} \geq P\{R \cdot \theta(X) \geq r\} \quad (100)$$

[**Todo:** voir figure a produire.]

Puisque u est concave^[Citation needed] et que B domine stochastiquement $R \cdot \theta(X)$, on a nécessairement que $\mathbf{E}u(sB) \geq \mathbf{E}u(R \cdot s\theta(X))$, pour tout $s > 0$. Or, par hypothèse de sous-linéarité on obtient que

$$\lim_{s \rightarrow \infty} \mathbf{E}u(R \cdot s\theta(X)) \leq \lim_{s \rightarrow \infty} u(sB) \quad (101)$$

$$= \lim_{s \rightarrow \infty} (\varrho u(-s\delta) + (1 - \varrho)u(s\bar{r}\xi)) \quad (102)$$

$$\leq \lim_{s \rightarrow \infty} -\varrho s\delta + (1 - \varrho)o(s) = -\infty, \quad (103)$$

ce qui démontre bien que s est borné.

5 Distributions synthétiques

6 Conclusion

SVM multiclasse

Time series et learning