

# Portfolio Optimization in a Big Data Context

Thierry BAZIER-MATTE

Summer 2015

## Abstract

Following [1], we provide a portfolio optimization method based on machine learning methods.

## 1 Introduction

This document considers a two-asset portfolio, of which one is the risk-free asset, yielding a constant return rate  $R_f$ , and the other being a risky asset  $s$ , typically a stock, yielding a random return rate  $r_{st}$  for each period  $t$ . The distribution underlying this return is assumed to be unknown. We suppose that each risky asset  $s$  can be described daily by an *information vector*  $x_{st}$  containing potentially useful information, such as technical, fundamental or news-related information. Furthermore, we assume that the allocation of each asset of the portfolio  $p_{st}$  can be fully determined using a *decision vector*  $q$  based on past information of the previous returns and of the previous information vectors. The allocation rule is the following:  $q^T x_{st}$  is allocated to the risky asset and  $1 - q^T x_{st}$  is allocated to the risk-free asset. Over the period  $t$ , the portfolio  $p_{st}$  consisting of asset  $s$  will therefore yield a return rate of:

$$p_{st}(q) = r_{st}q^T x_{st} + (1 - q^T x_{st})R_f. \quad (1)$$

The question we now wish to ask is how the decision vector  $q$  should be chosen. We assume we have access to a training dataset  $S_n$ , comprising of  $s$  different assets over  $t$  periods, such that  $n = s \times t$ .

What's the difference between having  $n$  points with having  $s \times t$  points? For example, what if  $s \gg t$  or the reverse?

## 2 Definitions and Bounds

### 2.1 Definitions Notation

Most of the following notation and definitions follow directly from [2].

Let  $S_n$  be a set of  $n$  vectors of  $\mathbf{R}^p \times \mathbf{R}$  of the form:

$$S_n = \{(x_1, r_1), \dots, (x_n, r_n)\}. \quad (2)$$

Each component of  $S_n$  is a tuple  $(x, r)$ , where  $x$  is the information vector and  $r$  is the observed return rate.

Using  $S_n$ , we wish to create a decision vector  $q_{S_n} \in \mathbf{R}^p$  from which we can make an investment decision  $p$  when confronted with a random draw  $d = (x, r)$ .

**Loss, cost and regret.** We introduce the loss  $\ell$  as being a function mapping from the space of the decision vectors and full observations (including market observations and the return of the asset) to a numerical quantity called the cost. Formally, given  $q$  a decision vector and a random draw  $d = (x, r)$ , the loss will be

$$\ell(q, d) = c(q(x), r) = c(q^T x, r). \quad (3)$$

Supposing an utility  $U$ , there are two different cost functions we can use. The first one, and perhaps the most obvious is defined by

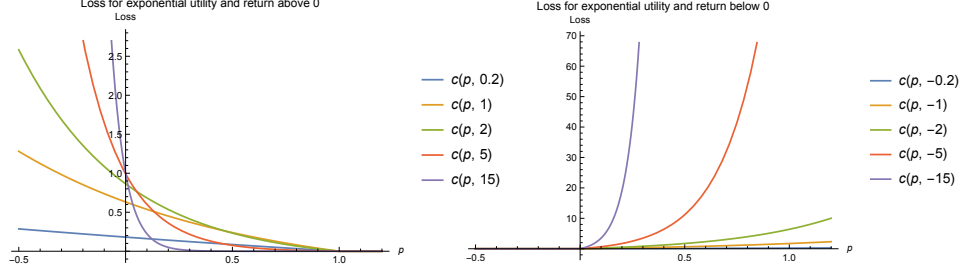
$$c(p, r) = -U(pr + (1 - p)R_f). \quad (4)$$

If the utility is piecewise linear, then the cost is not bounded, and so an infinite return would yield a negatively infinite cost. This means that risk-taking, ie. positions with  $|p| > 1$ , will be encouraged since they can in fact yield a negative cost. This also means that with an unbounded utility, the algorithm minimizing the cost over the sample must have a non-zero regularization term, otherwise the problem might be unconstrained. **Unless we add a constant piece if  $r < \bar{r}$ .**

We can also define another cost function that is always non-negative, that we shall call the regret  $\rho$  of the decision:

$$\rho(p, r) = \begin{cases} [U(r) - U(pr + (1 - p)R_f)] & \text{if } r > R_f \\ [U(R_f) - U(pr + (1 - p)R_f)] & \text{if } r \leq R_f, \end{cases} \quad (5)$$

where by  $[\cdot]$  we mean  $\max(\cdot, 0)$ . In the case of exponential utility, this cost function is actually equivalent to the one previously defined, but is quite different in the case of an unbounded utility function, for example the piece-wise linear one. Such a cost does not encourage risky position, since the cost is at least 0. This means that if  $r > 0$ , then there's no point having  $p > 1$  instead of  $p = 1$  since both will yield a zero cost position.



**Utility.** There are two ways we can model our utility, and both are concave function, to represent a risk-averse approach. The first utility is the linear utility of the form

$$U(r) = r + \min(0, \beta r), \quad (6)$$

with  $\beta > 0$ . The other utility is exponential:

$$U(r) = -\exp(-\mu r), \quad (7)$$

with  $\mu > 0$ .

**Algorithm.** We will be concerned with probabilistic confidence bounds on results produced using the following algorithm, using dataset  $S_n$ . We shall denote this algorithm as  $A_{S_n}$ .

$$q^* = \arg \min_{q \in \mathbf{R}^p} \frac{1}{n} \sum_{i=1}^n c(q^T x_i, r_i) + \lambda \|q\|_2^2. \quad (8)$$

**Assumptions.** We will assume that information vectors have been pre-processed and lie in a  $X_{\max}^2$  radius ball. We also assume that the return rates observed are comprised within  $[-\bar{r}, \bar{r}]$ . This last assumption will be relaxed.

## 2.2 Theoretical Optimal Value [Relevant?]

Before we go on to derive an algorithm for the in-sample optimal allocation, we first derive some theoretical results for an optimal allocation scalar  $p^*$ , provided the distribution of the returns is priorly known.

Formally, we suppose that this distribution is characterized by its CDF  $F$ . We then wish to minimize the expected cost  $c(p, r)$ , where  $p$  is a variable and  $r$  a random return rate. Since,

$$E[c(p, r)] = \int_{-\infty}^{\infty} c(p, r) dF(r), \quad (9)$$

the optimal allocation  $p^*$  is therefore the solution to

$$\frac{d}{dp} \int_{-\infty}^{\infty} c(p, r) dF(r) = 0. \quad (10)$$

However, our theoretical assumptions are fundamentally different from this one, since we consider a joint distribution  $\mathbf{D} = \mathbf{X} \times \mathbf{R}$ , so that returns are dependant of a draw from  $\mathbf{X}$ , whereas it's assumed here that decision can only be made on the basis of the returns.

## 2.3 Definitions, Theorems and Tools

We now state a certain number of definitions and theorems that will be of great importance in the analysis on the results produced by the algorithm mentioned above. These definitions and theorems are taken verbatim from [2].

**Definition.** An algorithm  $A$  has *uniform stability*  $\beta$  with respect to the loss function  $\ell$  if, for all  $S_n \in \mathbf{D}^n$  and  $i \in \{1, \dots, n\}$ , the following holds:

$$\|\ell(A_{S_n}, \cdot) - \ell(A_{S_n^{\setminus i}}, \cdot)\|_{\infty} \leq \beta_n, \quad (11)$$

or, equivalently,

$$\sup_{d \in \mathbf{D}} |\ell(A_{S_n}, d) - \ell(A_{S_n^{\setminus i}}, d)| \leq \beta_n. \quad (12)$$

Here,  $S_n^{\setminus i}$  means the set  $S_n$  with the  $i$ th data point removed.

Furthermore,  $A$  is *stable* when  $\beta_n = O(1/n)$ .

**Definition.** A loss function  $\ell$  is  $\sigma$ -*admissible* if the associated cost function  $c$  is convex with respect to its first argument and the following condition holds for any  $p_1, p_2$  and  $r$ :

$$|c(p_1, r) - c(p_2, r)| \leq \sigma |p_1 - p_2|. \quad (13)$$

It is easy to see that if  $\text{dom } c(\cdot, r)$  is bounded and  $c(\cdot, r)$  does not reach infinity, then its loss function is  $\sigma$ -admissible.

**Theorem 1.** Let  $F$  be a reproducing kernel Hilbert space with kernel  $\kappa$  such that  $\forall x \in X$ ,  $\kappa(x, x) \leq \kappa^2 < \infty$ . If  $\ell$  is  $\sigma$ -admissible with respect to  $F$ , then the learning algorithm defined by

$$A_S = \arg \min_{g \in F} \frac{1}{n} \sum_{i=1}^n \ell(g, d_i) + \lambda \|g\|_k^2 \quad (14)$$

has uniform stability  $\alpha_n$  with respect to  $\ell$  with

$$\alpha_n \leq \frac{\sigma^2 \kappa^2}{2\lambda n}. \quad (15)$$

**Definition.** The *true risk* with respect to algorithm  $A$  and set  $S_n$  is defined as

$$R_{\text{true}}(A, S_n) = E_d[\ell(A_{S_n}, d)], \quad (16)$$

which is, in plain words, the expected loss incurred when applying the algorithm created from training set  $S_n$  in the wild, ie. out of sample.

**Definition.** The *empirical risk* with respect to algorithm  $A$  and set  $S_n$  is defined as

$$\hat{R}(A, S_n) = \frac{1}{n} \sum_{i=1}^n \ell(A_{S_n}, d_i), \quad (17)$$

which is, in plain words, the average cost incurred by our model over all the training set.

**Theorem 2.** *Let  $A$  be an algorithm with uniform stability  $\alpha_n$  with respect to a loss function  $\ell$  such that  $0 \leq \ell(A_{S_n}, d) \leq M$  for all  $d = (x, r) \sim D$  and all sets  $S_n$  of size  $n$ . Then for any  $n \geq 1$  and any  $\delta \in (0, 1)$ , the following bound holds with probability at least  $1 - \delta$  over the random draw of the sample  $S_n$ :*

$$|R_{\text{true}}(A, S_n) - \hat{R}(A, S_n)| \leq 2\alpha_n + (4n\alpha_n + M) \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (18)$$

## 2.4 Remarks and applications to our algorithm

In this section we shall derive bounds on the proposed algorithms concerned with the difference between the in-sample average cost using the decision vector determined by our algorithm with the (theoretical) expected cost over the whole distribution when using the same decision vector.

The analysis will be done in two parts, one for each of our cost functions: first with  $c = -U$  (the canonical cost) and second with the cost function as the regret at each decision. Both are quite similar in their results.

### 2.4.1 Canonical cost.

We first consider the case where

$$c(p, r) = -U(pr + (1 - p)R_f). \quad (19)$$

**Remark** ( $\sigma$ -admissibility). We first remark that both forms of  $U$  yield a convex function of  $p$  with  $r$  fixed, since  $U$  is concave.

Next, the expression  $|c(p_1, r) - c(p_2, r)|$  reduces to

$$|U(p_1r + (1 - p_1)R_f) - U(p_2r + (1 - p_2)R_f)|. \quad (20)$$

Now because  $r \in [-\bar{r}, \bar{r}]$ ,  $U$  is Lipschitz continuous on its domain, and so (20) is bounded by

$$\alpha|p_1r + (1 - p_1)R_f - (p_2r + (1 - p_2)R_f)| = \alpha|p_1 - p_2||r - R_f| \quad (21)$$

where

$$\alpha = \sup_{r \in [-\bar{r}, \bar{r}]} |U'(r)|. \quad (22)$$

In the linear case, the derivative is piecewise constant, and is set to 1 on for returns below  $r_c$ , so that  $\alpha = 1$ . In the exponential case,  $U'(r) = \exp \mu r$ , and  $\alpha = \exp \mu \bar{r}$ .

The bound (21) must hold for any  $r$ . The expression  $|r - R_f|$  will reach its largest value at  $r = -\bar{r}$ , since  $R_f$  is assumed to be non-negative, and therefore

$$|c(p_1, r) - c(p_2, r)| \leq \alpha(\bar{r} + R_f)|p_1 - p_2|, \quad (23)$$

which shows that  $\ell$  is  $\sigma$ -admissible with the canonical cost, with  $\sigma = \alpha(\bar{r} + R_f)$ .

We also remark that the utility is actually invariant to scaling, ie.  $\tilde{U} = kU$  where  $k$  is positive yields the same relative utility. The  $\sigma$  bound can therefore be as tight as desired, since it can be scaled by an arbitrary positive constant, and the algorithm will yield the same result.

**Remark** (Algorithmic Stability). Because the loss function is  $\sigma$ -admissible, the conditions for Theorem 1 are respected: if we take the kernel on the space  $\mathbf{R}^p$  as the regular inner product, then  $\kappa(x, x) \leq X_{\max}^2$ . Therefore, our algorithm has algorithmic stability with:

$$\alpha_n \leq \frac{k^2(\bar{r} + R_f)^2 X_{\max}^2}{2\lambda n} \quad (24)$$

with linear utility and

$$\alpha_n \leq \frac{k^2 \exp(2\mu \bar{r})(\bar{r} + R_f)^2 X_{\max}^2}{2\lambda n} \quad (25)$$

in the case of exponential utility.

As discussed above,  $k > 0$  is a scaling term we can add to the utility function, with which we can actually have algorithmic stability as tight as needed.

How can we leverage  $k$ ?

**Remark** (Generalization Bound). Using Theorem 2, we can derive an out-of-sample bound on the result produced by our algorithm for each utility forms that holds with probability  $1 - \delta$ , for  $0 < \delta < 1$ . The  $M$  value stated by Theorem 2 corresponds to the highest loss we can incur on the domain, which happens when  $p = 1$  and  $r = -\bar{r}$ , ie.

$$M = c(1, -\bar{r}) = -U(\bar{r}). \quad (26)$$

In the linear utility case,

$$|R_{\text{true}}(A, S_n) - \hat{R}(A, S_n)| \leq 2\alpha_n + (4n\alpha_n + M) \sqrt{\frac{\log(2/\delta)}{2n}} \quad (27)$$

$$\leq \frac{k^2(\bar{r} + R_f)^2 X_{\max}^2}{\lambda n} + \left( \frac{2k^2(\bar{r} + R_f)^2 X_{\max}^2}{\lambda} + M \right) \sqrt{\frac{\log(2/\delta)}{2n}}, \quad (28)$$

and in the exponential utility case,

$$|R_{\text{true}}(A, S_n) - \hat{R}(A, S_n)| \leq 2\alpha_n + (4n\alpha_n + M)\sqrt{\frac{\log(2/\delta)}{2n}} \quad (29)$$

$$\leq \frac{k^2 \exp(2\mu\bar{r})(\bar{r} + R_f)^2 X_{\max}^2}{\lambda n} + \left( \frac{2k^2 \exp(2\mu\bar{r})(\bar{r} + R_f)^2 X_{\max}^2}{\lambda} + M \right) \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (30)$$

**Remark** (Absolute Bound). Finally, we derive a measure of performance when we compare  $\hat{q} = A_{S_n}$  with the true optimal  $q^*$  over the distribution, and how each perform when applied on the true risk of the distribution...

#### 2.4.2 Regret cost.

We now turn our attention toward the regret cost, ie.

$$c(p, r) = \lfloor U(\max(r, R_f)) - U(pr + (1 - p)R_f) \rfloor. \quad (31)$$

This type cost is different from the canonical cost in that it does not encourage ‘risky’ position for which  $0 < p < 1$ , since there’s no gain in over-investing our portfolio in a winning asset. Conversely, there’s no gain in shorting a losing asset. Such a cost could therefore be more suitable in a conservative environment, even if the results produced by the algorithm  $A$  might very well yield a position  $0 < p < 1$ .

**Remark** ( $\sigma$ -admissibility). The  $\sigma$ -admissibility of the regret cost follows closely from the  $\sigma$ -admissibility of the canonical cost derived above. Indeed, we remark that, provided that  $c(p_1, r), c(p_2, r) > 0$ , then the expression  $|c(p_1, r) - c(p_2, r)|$  reduces to

$$|U(p_1r + (1 - p_1)R_f) - U(p_2r + (1 - p_2)R_f)|, \quad (32)$$

which has exactly the same form as (20), and therefore yields the same  $\sigma$  bound. Now for the case where, without loss of generality  $c(p_2, r) = 0$ . Let  $\tilde{c}$  be defined as  $c$ , but without  $\lfloor \cdot \rfloor$ , so that  $\lfloor \tilde{c}(p, r) \rfloor = c(p, r)$ . In particular  $\tilde{c}(p_2, r) < 0$  and so

$$|c(p_1, r) - c(p_2, r)| = |\lfloor \tilde{c}(p_1, r) \rfloor - \lfloor \tilde{c}(p_2, r) \rfloor| \quad (33)$$

$$\leq |\tilde{c}(p_1, r) - \tilde{c}(p_2, r)| \quad (34)$$

$$= |U(p_1r + (1 - p_1)R_f) - U(p_2r + (1 - p_2)R_f)|, \quad (35)$$

which has again the same form as (20).

This means that the loss defined with the regret cost has the same  $\sigma$ -admissibility as the loss defined with the canonical cost, namely

$$\sigma = k(\bar{r} + R_f) \quad (36)$$

in the case of linear utility and

$$\sigma = k(\bar{r} + R_f) \exp(\mu \bar{r}) \quad (37)$$

in the case of exponential utility.

**Remark** (Algorithmic Stability). Because this regret cost has the same  $\sigma$ -admissibility as the canonical cost, then the algorithmic stability also remains the same for both algorithm, ie.

$$\alpha_n \leq \frac{k^2(\bar{r} + R_f)^2 X_{\max}^2}{2\lambda n} \quad (38)$$

with linear utility and

$$\alpha_n \leq \frac{k^2 \exp(2\mu \bar{r})(\bar{r} + R_f)^2 X_{\max}^2}{2\lambda n} \quad (39)$$

in the case of exponential utility.

**Remark** (Generalization Bound). Likewise, the only difference here with the canonical cost is the value of the  $M$  variable of Theorem 2, defined as the highest cost on the domain. In the case here, this value is again reached for an investment  $p = 1$  on a losing asset at  $r = -\bar{r}$ , so that,

$$M = \lfloor U(R_f) - U(\bar{r}) \rfloor. \quad (40)$$

We are left with the same expression for the generalization bounds, in the linear utility case,

$$|R_{\text{true}}(A, S_n) - \hat{R}(A, S_n)| \leq 2\alpha_n + (4n\alpha_n + M) \sqrt{\frac{\log(2/\delta)}{2n}} \quad (41)$$

$$\leq \frac{k^2(\bar{r} + R_f)^2 X_{\max}^2}{\lambda n} + \left( \frac{2k^2(\bar{r} + R_f)^2 X_{\max}^2}{\lambda} + M \right) \sqrt{\frac{\log(2/\delta)}{2n}}, \quad (42)$$

and in the exponential utility case,

$$|R_{\text{true}}(A, S_n) - \hat{R}(A, S_n)| \leq 2\alpha_n + (4n\alpha_n + M) \sqrt{\frac{\log(2/\delta)}{2n}} \quad (43)$$

$$\leq \frac{k^2 \exp(2\mu \bar{r})(\bar{r} + R_f)^2 X_{\max}^2}{\lambda n} + \left( \frac{2k^2 \exp(2\mu \bar{r})(\bar{r} + R_f)^2 X_{\max}^2}{\lambda} + M \right) \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (44)$$

## 2.5 True Optimal Bound

One bound that is still of interest is the absolute difference between the in-sample cost using in-sample decision  $\hat{R}(A, S_n)$  and the expected cost over the unknown distribution  $D$  (ie. the true risk) using the unknown optimal decision  $q^*$  defined by

$$q^* = \arg \min_{\{q: X \subset \mathbf{R}^p \rightarrow \mathbf{R}\}} R_{\text{true}}(q). \quad (45)$$



The true optimal bound is therefore expressed by

$$|R_{\text{true}}(q^*) - \hat{R}(A, S_n)|. \quad (46)$$

We can then apply the triangle inequality to bound this expression by the generalization bound plus  $|R_{\text{true}}(q^*) - R_{\text{true}}(\hat{q})|$ .

*However, whereas [1] could bound this expression using a theorem from quantile regression theory, I'm here at loss (no pun intended...)*

**Things to consider.** Here are some points that should be considered.

- We have assumed so far we were dealing with datasets taken iid. But in a financial setting, this is not necessarily the case, as influence of features might be changing with time. We would then have independent, but not necessarily identical draws. We could (and should) however assume that underlying probability distribution is changing *smoothly* (moment-generating functions, Kalman filters, change detection, etc.)
- In the newsvendor article, the generalization bounds actually hold with probability  $1 - \delta - n\gamma$ , if we assume  $\bar{D}$  (the highest considered demand) holds with probability  $1 - \gamma$ . It is only mentioned in the non-regularized case ([1], p. 14). A reference would be needed, as the  $n\gamma$  term grows up to be of no use [Reference needed !]
- We are also considering discrete quantities, updated daily. Our model should encompass quantities that might change during the day.
- A section on how data should be normalized must also be added. This is essential to go from theory to practical applications.
- Not to forget: the  $p$  features should also include higher order terms of the raw features, ie.  $x_i^2$  and  $x_i x_j$  so we can capture a second order Taylor approximation (see [1], p. 6). Instead of  $p$  features, we would then have  $p^2 + p$  features.
- The fact we are constraining features to live within a  $X_{\text{max}}$  radius, ie.  $\|x\|_2 \leq X_{\text{max}}$ , is something that we would like to relax. (Like changing the norm, or perhaps the metric?)

### 3 Multi-asset Portfolio

We now consider a multi-asset portfolio, made of  $m$  risky assets and a risk-free asset. Each of these risky assets have information vector  $x_{st}$ , where  $s \in \{1, \dots, m\}$  and  $t \in \{1, \dots, n\}$ . Therefore, we suppose having access to a sample dataset  $S_n = \{X_1, \dots, X_n\}$  where  $X_i \in \mathbf{R}^{m \times p}$ .

There are two ways to approach the multi-asset scenario.

#### 3.1 First approach

On one hand, we can define a single decision vector  $q$  such that the return on a single holding period would be defined by

$$r^T Xq + (1 - 1^T Xq)R_f, \quad (47)$$

where  $r \in \mathbf{R}^m$  is the return vector of the  $m$  stocks and  $1^T$  is a constant vector of ones in  $\mathbf{R}^m$  whose role is simply to sum over the components of  $Xq$ , ie. sum over the risky allocations of the portfolio. However, this method has the disadvantage that the decision vector is applied for each stock in the same way, and so all their information vectors are ultimately ‘averaged’ over. Therefore, if some asset is more sensitive to an information component than other assets, this knowledge might be lost if using this method.

Nevertheless, we can again define two different cost functions, which our optimal decision vector would minimize when applied to our sample dataset  $S_n$ . Again, the cost can be simply defined as

$$c(p, r) = -U(r^T p + (1 - 1^T p)R_f). \quad (48)$$

As discussed previously, this definition might have the disadvantage of providing an unconstrained optimization problem when determining the optimal decision vector.

We are not limited to a single cost definition, and in fact, the multi-asset scenario provides many legitimate cost definitions. Here the most ‘aggressive’ one:

$$c(p, r) = \begin{cases} [U(\max r) - U(r^T p + (1 - 1^T p)R_f)] & \text{if } \max r > R_f \\ [U(R_f) - U(r^T p + (1 - 1^T p)R_f)] & \text{if } \max r \leq R_f. \end{cases} \quad (49)$$

This means that we expect (or conversely, that we suffer a lesser loss) when the total return of our portfolio is near the best return achieved by our considered stocks. Instead of  $\max r$ , we also could have chosen the average returns  $\bar{r}$  during the period as the target that we wish to achieve (provided that the average is superior to the risk-free return). We could in fact choose any function mapping from a vector of returns  $r$  to a scalar as our global return objective.

In any case, the algorithm determining the decision vector  $q$  would once be defined as previously, ie. as (8).

### 3.2 Second approach

Under this second approach, instead of averaging over all information vector at the risk of losing asset-specific information, we now use a decision matrix  $Q$  of size  $p \times m$ , whose  $m$  columns are actually trained independently one to each other so that each asset reacts differently from the rest.

The allocation of each vector would then be  $\text{diag}(XQ)$ , ie. the diagonal elements of the  $m \times m$  matrix  $XQ$ . The one-period return is therefore

$$r^T \text{diag}(XQ) + (1 - 1^T \text{diag}(XQ))R_f. \quad (50)$$

This method still has a disadvantage though, in that it won't consider the eventual influence of information  $x_j$  of asset  $j$  to asset  $j$ , ie. each stock is trained independantly one to each other.

### 3.3 Third Approach

Finally, we introduce a third approach with the objective of using all available data of the market, available through  $X \in \mathbf{R}^{m \times sp}$ . However, we must be clear about what information  $X$  carries with it; its information is only bounded to different assets, and therefore does not include 'global' information with it.

Following a simple idea introduced by [5], we can add an unknown linear transformation  $W$  to  $X$  so that  $q$  can benefit from the whole information matrix, and not only from the average of its columns. The total return using a transformation  $W$  and decision vector  $q$  would therefore be expressed as

$$r^T XWq + (1 - 1^T XWq)R_f. \quad (51)$$

### 3.4 Fourth Approach

We now introduce a more general framework that actually encompasses all other methods we've seen so far. In addition, we also add regularization terms so that the following rules are respected :

- Given an asset  $i$ , we insist that the allocation scalar for  $i$  relies more strongly on feature vector  $x_i$  than on feature vectors from other assets in the market.
- The market should be segmented in different industries, and our algorithm should partition the  $m$  stocks in  $k$  groups of similar nature. This would also impact how cross-influence of feature affect very weakly assets of other groups.
- Each asset  $i$  has its own decision vectors  $q_{ij}$ , where  $q_{ij}$  is the influence of asset  $j$  on asset  $i$ .

With this in mind we can now define a general allocation rule given feature vectors  $\{x_i\}$  for each asset in the market. As before, we consider a general inner product  $x_i^T q_j$  as a mapping from the feature space of  $i$  to the allocation scalar space. Therefore, the total influence of the  $m$  assets on  $i$  can be expressed as

$$x_i^T q_{ii} + \sum_{j \neq i} x_j^T q_{ij}. \quad (52)$$

We can reduce this expression by introducing a decision matrix  $Q_i \in \mathbf{R}^{p \times m}$  defined by

$$Q_i = ([Q_i]_1 \quad \cdots \quad [Q_i]_m). \quad (53)$$

By using the same feature matrix  $X$ , the above expression is therefore equivalent to

$$\text{tr}(XQ_i), \quad (54)$$

so that the total return will be given by

$$\sum_{i=1}^m r_i \text{tr}(XQ_i) + \left(1 - \sum_{i=1}^m \text{tr}(XQ_i)\right) R_f. \quad (55)$$

Now, as before the  $Q_i$  matrices are unknown at first, but can be learned using an optimization problem, where the objective is the average cost  $\bar{c}$ , defined by (48). Let us, just for a moment, neglect the regularization term. Then, the optimal matrices are  $Q_i^*$  are the solution to the following optimization problem:

$$\begin{aligned} & \underset{Q_1^*, \dots, Q_m^*}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n c(p_i, r_i) \\ & \text{s. t.} \quad [p_i]_j = \text{tr}(X_i Q_j). \end{aligned} \quad (56)$$

Now, there are some fundamental questions concerning this optimization problem which we'll report for later. Indeed, before we start analyzing the theoretical guarantees we can get out of this problem, we would first like to add regularization terms to the objective, with the hope of being compliant with the rules we've proposed above.

The first rule we proposed is certainly the most easy to implement. What it says is basically that in the allocation scalar for asset  $i$ , the contribution from feature  $i$  should be larger than contribution coming from feature of other assets.

Mathematically, this would translate in this 'soft' constraint:

$$x_i^T [Q_i]_i \gg \sum_{j \neq i} x_j^T [Q_i]_j. \quad (57)$$

In an optimization context, we can add the following expression to the objective:

$$+ \text{tr}(XQ_i)^2 - \lambda (x_i^T [Q_i]_i)^2, \quad (58)$$

with  $\lambda \geq 0$ . In plain words, this means that we expect the whole allocation scalar close to the contribution from  $x_i^T[Q_i]_i$ , with a regularization parameter  $\lambda$  to avoid all other terms  $x_j^T[Q_i]_j$  falling to zero.

Now, let's tackle the second constraint, where we expect the cross-contribution of a feature vector to only affect a close number of assets. For example, we'd expect technology assets to affect each other, but with little impact on energy assets. This evidently reminds of  $L1$  regularization, where the activation is non-smooth, and either you're part of a group, or you're not.

We can first start by defining the number of groups, or segments, we expect the market to have. We shall note this number of groups by  $k$ .

Let us then state the problem formally. Let  $\Xi \in \mathbf{R}^{m \times m}$  be the allocation matrix, such that

$$[\Xi]_{ij} = x_j^T[Q_i]_j. \quad (59)$$

Then, the allocation vector  $p$  for each of the  $m$  assets corresponds to the sum of each of the rows of  $\Xi$ . Therefore,  $\Xi$  is invariant to reordering of individual rows. In the process of optimization, we want to have to turn  $\Xi$  in a  $k$ -block matrix, with each block representing a certain sector. However, this *block* requirement should not be too hard, as in certain event, we could have assets with overlapping sectors. Anyway,  $L1$  sounds like a good idea, since this kind of regularization often allows for features to be either 'in' or 'out'. In regularization, there's always two opposing forces, and its role is to find balance between those two forces.

So what are the two forces at play here? Let us consider for one moment a certain asset  $i$ . Then, its allocation  $p_i$ , given a feature matrix  $X$  whose rows can be expressed by  $x_j^T$ , will be of the form

$$p_i = \text{tr}(XQ_i) = \sum_{j=1}^m x_j^T[Q_i]_j. \quad (60)$$

The optimization algorithm, in order to find the optimal  $Q_i$ , no matter what the shape of the utility is, will try to increase  $p_i$  as much as possible if  $r_i > 0$ , and conversely, will try to decrease it if  $r_i < 0$ . Let's pretend for a moment that  $r_i > 0$ , so that we're trying to drive  $p_i$  as positive as possible.

★ ★ ★

With  $\Xi$  defined as above, than we can state the optimal  $\Xi$  as an optimization problem of the following form:

$$\text{minimize} \quad -U(1^T \Xi r), \quad (61)$$

where  $\Xi$  is a function of  $m^2$  independent vectors  $q_1^1, \dots, q_m^1, \dots, q_1^m, \dots, q_m^m$ , such that  $\Xi$  can be expressed by

$$\Xi = \begin{pmatrix} x_1^T q_1^1 & x_2^T q_2^1 & \cdots & x_m^T q_m^1 \\ x_1^T q_1^2 & x_2^T q_2^2 & \cdots & x_m^T q_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^T q_1^m & x_2^T q_2^m & \cdots & x_m^T q_m^m \end{pmatrix}. \quad (62)$$

We can further decompose this expression. First we need to introduce a decision matrix  $Q \in \mathbf{R}^{m \times mp}$ :

$$Q = \begin{pmatrix} q_1^{1T} & q_2^{1T} & \cdots & q_m^{1T} \\ q_1^{2T} & q_2^{2T} & \cdots & q_m^{2T} \\ \vdots & \vdots & \ddots & \vdots \\ q_1^{mT} & q_2^{mT} & \cdots & q_m^{mT} \end{pmatrix}, \quad (63)$$

as well as a feature matrix  $X$ :

$$X = \begin{pmatrix} x_1 & 0 & 0 & 0 \\ 0 & x_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & x_m \end{pmatrix} \quad (64)$$

where each  $x_i \in \mathbf{R}^p$  should be understood as column vectors, so that  $X \in \mathbf{R}^{mp \times m}$ .

With these definitions, we have

$$\Xi = QX. \quad (65)$$

The allocation vector  $p$  can also be expressed as

$$p^T = 1^T \Xi = 1^T QX, \quad (66)$$

so that the total return with asset returns vector  $r$  is

$$1^T QXr. \quad (67)$$

This is the quantity we seek to maximize, or rather with which we seek to cost-maximize (More on that later).

★ ★ ★

Now what we want to do is to solve  $Q$ . For now, we can try to simply optimize  $1^T QXr$ , with  $Q$  being the free variable.

This problem is not bounded, for we can simply target positive  $r_i$ , and assign them positive (infinite)  $p_i$  weight, and conversely for negative  $r_i$  to obtain an infinite return. Therefore, we want to regularize the weight from  $Q$ , and we can do so by adding a regularization of the form

$$+ \lambda \|\Xi\| \quad (68)$$

for some appropriate matrix norm  $\|\cdot\|$ .

Now, concerning the block structure of  $\Xi = QX$ . We shall try to solve this problem iteratively. We first suppose that  $\Xi$  is already of the correct form, ie. the columns have the right order, but we don't know yet how they should be divided in blocks.

The idea is to define a recursive function  $\rho_\Xi$  of  $\Xi$  defined like so:

$$\rho_\Xi(1) = \Xi_{11} \quad (69)$$

$$\rho_\Xi(2) = \begin{pmatrix} \rho_\Xi(1) & \Xi_{12} \\ \Xi_{21} & \Xi_{22} \end{pmatrix} \quad (70)$$

$$\dots = \dots \quad (71)$$

This means that  $\rho_\Xi(k)$  represents the first  $k$ -block of  $\Xi$ , but defined in term of previous  $\rho_\Xi$  definitions.

With such a recursive function, we can introduce a progressive regularization parameter that grows with each new block, to represent that adding a new row/column to the block should be 'harder' at each iteration. We would therefore have a total regularization term of the form

$$+ \sum_{i=1}^m f(i) \|\rho_\Xi(i)\|_1, \quad (72)$$

where  $f$  is a wisely chosen monotonic and positive function that depends on  $k$ , the number of blocks ( $k$  is a parameter of the problem, ie. we tell how many blocks we expect). One constraint this add is that  $f(i) = \infty$  for values  $i > m - k$ . For values under  $m - k$ ,  $f(i)$  could be of the form

$$f(i) = \lambda^i. \quad (73)$$

## **I might have a rough idea...**

*Following is a collection of ideas. This is probably to be taken not very seriously.*

Add something perhaps on how he can cramp into the model some parameters so that only a handful assets are being chosen (perhaps L1 reg. or something like that). Also maybe we could add another complexity where we would like to keep stocks we already have, perhaps using a utility telling how much we want to discourage the new daily assets. That sort of things.

Something else that was so far overlooked is the fact that the sum of all allocation components need not sum up to 1. Again, a regularization with a plateau shape form could perhaps translate what we need.



## References

- [1] Cynthia Rudin and Gah-Yi Vahn. *The Big Data Newsvendor: Pratical Insights from Machine Learning*, Operations Research, 2015.
- [2] Olivier Bousquet and André Elisseeff. *Stability and Generalization*, Journal of Machine Learning Research, 2002.
- [3] Stephen Boyd, Lieven Vandenberghe. *Convex Optimization*, Cambridge University Press, 2004.
- [4] Rockafellar, R. T. *Convex Analysis*, Princeton University Press, 1970.
- [5] Ming Fai Wong et al. *Stock Market Prediction from WSJ: Text Mining via Sparse Matrix Factorization*. ICDM 2014.