

The Big Data Newsvendor Problem in a Portfolio Optimization 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 . 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 . 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 $n = s \times t$.

2 Definitions and Bounds

2.1 Definitions and Notation

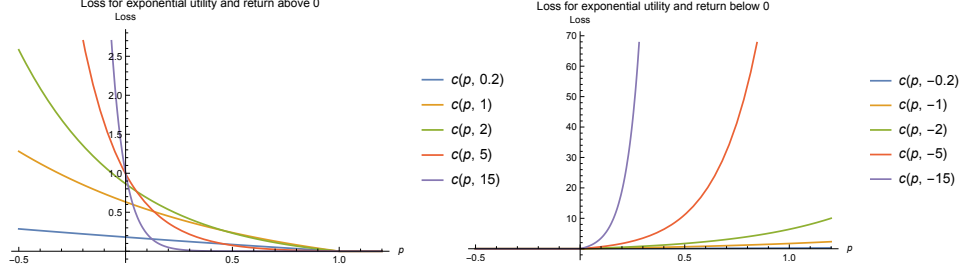
Most of the following notation and definitions follow directly from [2].

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

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

Maybe considerations about the length of the period should be added? For example, it's not specified what's the period length of R_f .

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



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 when confronted with a random draw $d = (x, r)$.

Loss and Cost. We introduce the loss ℓ and the cost c of using q with a random draw $d = (x, r)$:

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

The cost must always be a non-negative quantity. Supposing an utility U , we model it as follows:

$$c(p, r) = \begin{cases} \lfloor U(r) - U(pr + (1-p)R_f) \rfloor & \text{if } r > R_f \\ \lfloor U(R_f) - U(pr + (1-p)R_f) \rfloor & \text{if } r \leq R_f \end{cases} \quad (4)$$

By $\lfloor \cdot \rfloor$ we mean a function returning its argument if non-negative and zero otherwise. This means that we don't want to discourage taking risk (borrowing or short-selling), but it's not encouraged either.

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

$$U(r) = r + \min(0, \beta r), \quad (5)$$

with $0 < \beta < 1$. The other utility is exponential:

$$U(r) = -\exp(-\mu r), \quad (6)$$

with $\mu > 0$.

Algorithm. We will be concerned with probabilistic confidence bounds on results produced using the following algorithm, using dataset S_n .

$$q^* = \arg \min_{q \in \mathbf{R}^p} \frac{1}{n} \sum_{i=1}^n c(q^T x_i, r_i) + \lambda \|q\|_2^2. \quad (7)$$

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.

*Include
reference for
definitions
and theorems*

Definition. An algorithm A has uniform stability β with respect to the loss function ℓ if, for all $S \in \mathcal{D}^n$ and $i \in \{1, \dots, n\}$, the following holds:

$$\|\ell(A_S, \cdot) - \ell(A_{S \setminus i}, \cdot)\|_\infty \leq \beta_n. \quad (8)$$

Here, $S \setminus i$ means the set S with the i th data point removed.

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

Definition. A loss function ℓ is σ -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 :

$$|c(p_1, r) - c(p_2, r)| \leq \sigma |p_1 - p_2| \quad (9)$$

Remark. Our loss function ℓ is σ -admissible with $\sigma = \bar{r} + R_f$ in the linear case and $\sigma = (\bar{r} + R_f) \exp(\mu \bar{r})$ in the exponential case.

Proof. First, we remark that both forms of U yield a convex function of p with r fixed.

Now we'll suppose that $c(p_1, r), c(p_2, r) > 0$. Then the expression $|c(p_1, r) - c(p_2, r)|$ reduces to

$$|U(p_1 r + (1 - p_1)R_f) - U(p_2 r + (1 - p_2)R_f)|. \quad (10)$$

Now because $r \in [-\bar{r}, \bar{r}]$, U is Lipschitz continuous on its domain, and so (10) is bounded by

$$\alpha |p_1 r + (1 - p_1)R_f - (p_2 r + (1 - p_2)R_f)| = \alpha |p_1 - p_2| |r - R_f| \quad (11)$$

where

$$\alpha = \sup_{r \in [-\bar{r}, \bar{r}]} |U'(r)|. \quad (12)$$

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 (11) 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.

Finally we consider the case where, without loss of generality, $c(p_2, r) = 0$. Then, if c had not been defined using $\lfloor \cdot \rfloor$, then we would have

$$\begin{aligned} |\lfloor c(p_1, r) \rfloor - \lfloor c(p_2, r) \rfloor| &\leq |c(p_1, r) - c(p_2, r)| \\ &\leq \sigma |p_1 - p_2|. \end{aligned} \quad (13) \quad \square$$

Theorem 1. Let F be a reproducing kernel Hilbert space with kernel κ that $\forall x \in X$, $\kappa(x, x) \leq \kappa^2 < \infty$. If ℓ is σ -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 α_n with respect to ℓ with

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

Remark. Our proposed algorithm has the form (14), and so has algorithmic stability bounded by

$$\alpha_n \leq \frac{(\bar{r} + R_f)^2 X_{\max}^2}{2\lambda n} \quad (16)$$

with linear utility and

$$\alpha_n \leq \frac{\exp(2\mu\bar{r}) X_{\max}^2}{2\lambda n} \quad (17)$$

in the case of exponential utility.

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

$$R_{\text{true}}(A, S_n) = E_d[\ell(A_{S_n}, d)], \quad (18)$$

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

$$\hat{R}(A, S_n) = \frac{1}{n} \sum_{i=1}^n \ell(A_{S_n}, d_i), \quad (19)$$

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

Remark. The maximum loss we can suffer over a single data point happens when $r_i = -\bar{r}$ and $p = 1$, ie.

$$c(1, -\bar{r}) = U(R_f) - U(\bar{r}). \quad (20)$$

We shall call this quantity γ .

Theorem 2. Let A be an algorithm with uniform stability α_n with respect to a loss function ℓ 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 :

$$|R_{\text{true}}(A, S_n) - \hat{R}(A, S_n)| \leq 2\alpha_n + (4n\alpha_n + M) \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (21)$$

Remark. Our algorithm has a generalization bound of

$$|R_{\text{true}}|(A, S_n) - \hat{R}(A, S_n) \leq 2\alpha_n + (4n\alpha_n + \gamma) \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (22)$$

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] “Si la valeur absolue de la dérivée est majorée par k , f est k -lipschitzienne”. Application lipschitzienne.
- [4] Rockafellar, R. T. *Convex Analysis*, Princeton University Press, 1970.
- [5] Supergradients.
- [6] Reference needed!