

# Distributionally Robust Mean-Variance Portfolio Selection with Wasserstein Distances

Jose Blanchet\* Lin Chen† Xun Yu Zhou‡

July 31, 2021

## Abstract

We revisit Markowitz’s mean-variance portfolio selection model by considering a distributionally robust version, where the region of distributional uncertainty is around the empirical measure and the discrepancy between probability measures is dictated by the Wasserstein distance. We reduce this problem into an empirical variance minimization problem with an additional regularization term. Moreover, we extend the recently developed inference methodology to our setting in order to select the size of the distributional uncertainty as well as the associated robust target return rate in a data-driven way. Finally, we report extensive backtesting results on S&P 500 that compare the performance of our model with those of several well-known models including the Fama–French model and the Black–Litterman model.

**Key Words.** Mean-variance portfolio selection, robust model, Wasserstein distance, Robust Wasserstein Profile Inference.

## 1 Introduction

We study data-driven mean–variance portfolio selection with model uncertainty (or ambiguity). The classical one-period Markowitz mean–variance model (Markowitz (1952)) is to choose a portfolio weighting vector  $\phi \in \mathbb{R}^d$  (all the vectors in this paper are by convention columns) among  $d$  stocks to maximize the risk-adjusted expected return. The precise formulation is<sup>1</sup>

$$\min_{\phi \in \mathbb{R}^d} \left\{ \phi^\top \text{Var}_{\mathbb{P}^*} (R) \phi : \phi^\top \mathbf{1} = 1, \phi^\top \mathbb{E}_{\mathbb{P}^*} (R) = \rho \right\}, \quad (1)$$

---

\*Department of Management Science and Engineering, Stanford University, Stanford, California 94305, USA, [jblanche@stanford.edu](mailto:jblanche@stanford.edu). Material in this paper is based upon work supported by the Air Force Office of Scientific Research under award number FA9550-20-1-0397. Additional support is gratefully acknowledged from NSF grants 1915967, 1820942, and 1838576.

†Department of Industrial Engineering and Operations Research, Columbia University, New York, New York 10027, USA, [lc3110@columbia.edu](mailto:lc3110@columbia.edu)

‡Department of Industrial Engineering and Operations Research, Columbia University, New York, New York 10027, USA, [xz2574@columbia.edu](mailto:xz2574@columbia.edu)

<sup>1</sup>There are several mathematically equivalent formulations of the original mean–variance model.

where  $R$  is the  $d$ -dimensional vector of random returns of the stocks;  $\mathbb{P}^*$  is the probability measure underlying the distribution of  $R$ ;  $\mathbb{E}_{\mathbb{P}^*}$  and  $\text{Var}_{\mathbb{P}^*}$  are respectively the expectation and variance under  $\mathbb{P}^*$ ;  $\rho$  is the targeted expected return of the portfolio.

It is well known that this model has a major drawback when applied in practice. On one hand, its solutions are very sensitive to the underlying parameters, namely the mean and the covariance matrix of the stocks. On the other hand,  $\mathbb{E}_{\mathbb{P}^*}$  is unknown in practice; so one has to resort to the empirical versions of the mean and the covariance matrix instead, which are usually significantly deviated from the true ones (especially the mean, due to the notorious “mean-blur” problem).

This motivates the development of the “robust” formulation of the Markowitz model which recognizes and tries to account for the impact of the (potentially significant) discrepancies between  $\mathbb{P}^*$  and its empirical version. This idea originates in the robust control approach from control theory (see, for example, Petersen et al. (2000)). Hansen and Sargent (2008) give a systematic account on applications of robust control to economic models. There is also a rich literature on robustification of portfolio choice. The paper by Lobo and Boyd (2000) is among the first to provide a worst-case robust analysis with respect to the second-order moment uncertainty within the Markowitz framework. Goldfarb and Iyengar (2003) consider a robust Markowitz problem with the uncertainty set based on vector/matrix distance. Pflug and Wozabal (2007) present a Markowitz model with distributional robustness based on a Wasserstein distance, a metric measuring the discrepancy between two probability measures which we also apply in this paper.<sup>2</sup> Nevertheless, their formulation involves an additional value-at-risk type of constraint that leads to a much more complex optimization problem. More importantly, their choice of the uncertainty size is exogenous and no guidance for optimally selecting the size is given.

Esfahani and Kuhn (2018) provide representations for the worst-case expectations in a Wasserstein-based ambiguity set centered at the empirical measure, and then apply their results to portfolio selection using different risk measures, leading to models different from the Markowitz model. An important difference in our approach, which is related to the work in Blanchet et al. (2016), as we shall discuss, is that we focus on the order-2 Wasserstein distance. This is important because, due to the quadratic nature of the variance objective that we consider, applying an uncertainty set based on Wasserstein of order 1 could potentially lead to arbitrarily large variances. We also note that the work in Esfahani and Kuhn (2018) provides guidance to choose the uncertainty size,  $\delta$ . But this choice of the uncertainty size deteriorates substantially with an increase in the dimension of the underlying portfolio. So, as we shall elaborate, we employ an approach similar to that proposed in Blanchet et al. (2016), which must be adapted and extended to our setting. Our current work relates to the broad literature on distributionally robust optimization (DRO). In addition to Esfahani and Kuhn (2018), related duality results for Wasserstein DRO formulations in which the probability model appears linearly in the objective function have been studied in Zhao and Guan (2018), Esfahani and

---

<sup>2</sup>A Wasserstein distance is the optimal value for a specific optimal transport problem. The notion was first formulated by Monge (1781) and its theory developed by Kantorovich (1942). It has been widely used in the study of distributionally robust optimization (DRO) problems.

Kuhn (2018) and Gao and Kleywegt (2016). A general result (with conditions which match the standard assumptions of the general optimal transport theory) is given in Blanchet and Murthy (2019). These results are not directly applicable to our setting because our objective function is not linear, but quadratic in the underlying probability model, so the techniques need to be adapted.

Also in connection to robust portfolio optimization, we mention the work in Delage and Ye (2010), which constructs uncertainty regions involving means and covariances of the return vector. The paper of Wozabal (2012) also considers a robust portfolio model with risk constraints based on expected short-fall, resulting in an optimization problem that requires solving multiple convex problems. These papers do not consider an optimal choice of the size of the uncertainty sets, as we do here.

Then, there are a number of papers that address a wide range of optimization techniques (such as interior-point methods, conic programming and linear matrix inequalities) in solving robust portfolio selection problems. A selection of these include Halldorsson and Tutuncu (2003), Costa and Paiva (2002) and Ghaoui et al. (2003). We obtain formulations that can be solved with basically any standard convex optimization software. Finally, we mention the work of Goh and Sim (2010) and Wiesemann et al. (2014) who investigate different forms of distributional ambiguity sets, as well as those of Hu and Hong (2013) and Jiang and Guan (2016) who study distributional robust formulations based on the Kullback-Leibler divergences. It is worth noting that the Kullback-Leibler divergence-based formulation is popular in economics (see Hansen and Sargent (2008)). Our formulation focuses on the use of Wasserstein ambiguity sets because of the intuitive out-of-sample exploration induced by the Wasserstein distance and because of the regularization interpretation which, as we shall see, results from our formulation.

Precisely, in this paper, we are interested in studying a distributionally robust mean–variance (DRMV) model, given by

$$\min_{\phi \in \mathcal{F}_{\delta, \bar{\alpha}}(n)} \max_{\mathbb{P} \in \mathcal{U}_{\delta}(\mathbb{P}_n)} \left\{ \phi^{\top} \text{Var}_{\mathbb{P}}(R) \phi \right\}, \quad (2)$$

where  $\mathbb{P}_n$  is the empirical probability derived from historical information of the sample size  $n$ ,  $\mathcal{U}_{\delta}(\mathbb{P}_n) := \{\mathbb{P} : D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta\}$  is the ambiguity set,

$$\mathcal{F}_{\delta, \bar{\alpha}}(n) = \left\{ \phi : \phi^{\top} \mathbf{1} = 1, \min_{\mathbb{P} \in \mathcal{U}_{\delta}(\mathbb{P}_n)} [\mathbb{E}_{\mathbb{P}}(\phi^{\top} R)] \geq \bar{\alpha} \right\},$$

is the feasible region of portfolios,  $\mathbb{E}_{\mathbb{P}}$  and  $\text{Var}_{\mathbb{P}}(R)$  denote respectively the mean and the covariance matrix under  $\mathbb{P}$ , and  $D_c(\cdot, \cdot)$  is a notion of discrepancy between two probability measures based on a suitably defined Wasserstein distance.<sup>3</sup>

Intuitively, formulation (2) introduces an artificial adversary  $\mathbb{P}$  (whose problem is that of the inner maximization) as a tool to account for the impact of the model uncertainty around the empirical distribution. There are two key parameters,  $\delta$  and  $\bar{\alpha}$ , in this formulation, and

---

<sup>3</sup>Recent work by Blanchet et al. (2016) shows that a similar definition of discrepancy in some other models recovers exactly certain well-known machine learning algorithms, such as square-root Lasso and support vector machines.

they need to be carefully chosen. The parameter  $\delta$  can be interpreted as the power given to the adversary: The larger the value of  $\delta$  the more power is given. If  $\delta$  is too large relative to the evidence (i.e. the size of  $n$ ), then the portfolio selection will tend to be unnecessarily conservative. On the other hand,  $\bar{\alpha}$  can be regarded as the lowest acceptable target return given the ambiguity set. Naturally, the choice of  $\bar{\alpha}$  should be based on the original target  $\rho$  given in (1); but one also needs to take into account the size of the distributional uncertainty,  $\delta$ . Using  $\bar{\alpha} = \rho$  will tend to generate portfolios that are too aggressive; it is more sensible to choose  $\bar{\alpha} < \rho$  in a way such that  $\rho - \bar{\alpha}$  is naturally informed by  $\delta$ .

This paper makes three main contributions. First, we show that (2) is equivalent to an (explicitly formulated) non-robust minimization problem in terms of the empirical probability measure in which a proper penalty term or “regularization term” is added to the objective function. The explicit regularization term that is derived from (2) is given in Theorem 1 below. This connects (and contrasts) to the *directly introduced* use of regularization in variance minimization techniques widely employed both in the statistics/machine learning literature and in practice to, among others, address the issue of overfitting. Indeed, practitioners who use mean–variance portfolio selection models often introduce regularization penalties, inspired by Lasso, in order to enhance the sparsity so as to include fewer stocks in their portfolios. Our use of Wasserstein distance to model distributional uncertainty naturally gives rise to a regularization term, suggesting an alternative, yet theoretical, justification and interpretation for its use in practice. Our result shows that our robust strategies are able to enhance out-of-sample performance with basically the same level of computational tractability as the standard mean-variance selection.

Our second main contribution provides guidance on the choice of the size of the ambiguity set,  $\delta$ , as well as that of the worst mean return target,  $\bar{\alpha}$ . This is accomplished by adapting and extending the robust Wasserstein profile inference (RWPI) framework, recently introduced and developed by Blanchet et al. (2016), to our setting in a data-driven way that combines optimization principles and basic statistical theory under suitable mixing conditions on historical data.

The last contribution empirically compares the performance of our DRMV strategies with those of several well-known and well-practiced models including the classical Markowitz model, the Fama–French model and the Black–Litterman model. We also compare our strategies with those of another robust model, the one put forward by Goldfarb and Iyengar (2003) in which robustness is based on vector/matrix distance. All these models (including ours) are static, single-period ones whereas in practice, a stock market is highly dynamic. In our empirical experiments, we implement them in the same rolling horizon fashion to account for the market dynamics. It should be noted that our theory applies to only a one-period model, and the numerical implementation to the multi-period market is heuristic based on rolling horizons. Finally, we also include in our comparison pre-committed optimal strategies based on a well-calibrated, non-robust continuous-time model in Cui et al. (2012). The experiments are carried out on S&P 500 for the backtesting period 2000-2017 with the prior 10 years as the training period. Our experiments show that DRMV compares favorably against all other models in achieving no worse (far better against most models) average returns and much lower variabilities.

This, we believe, is another important insight that we can draw from this research and that merits further investigation.

We should acknowledge, however, that our approach fails when the number of stocks  $d$  is not small compared with the sample size  $n$ , since the inference method for choosing  $\delta$  and  $\bar{\alpha}$  theoretically relies on some central limit theorems which require  $n$  to asymptotically approach infinity. Moreover, when  $d > n$ , both the empirical probability measure and the plug-in portfolios are not consistent. The application in our mind is asset allocation in which relatively small portfolios (i.e. those having dozens of assets) are desired. Indeed, one of the main goals of the regularization technique is to reduce the number of stocks in a portfolio. In our empirical study we have 100 stocks and only 108 data points, and our algorithm achieves good performance.

Ao et al. (2018) also provide a theoretical justification on the  $L^1$ -norm regularization by proving that the regularized portfolios asymptotically achieve the optimal mean and variance when the size of the portfolio approaches infinity. However, their result relies on the normality of return distribution, whereas we only require the data to be stationary which in particular allows heavy tail distributions commonly seen in financial data. Ao et al. (2018) estimate the regularization coefficient –  $\delta$  in our setting – using the cross validation technique, while we use the statistical inference to derive this parameter in a data-driven way.<sup>4</sup> Moreover, all the results in Ao et al. (2018) seems to work only with the  $L^1$ -norm, whereas our results can be applied to any  $L^p$ -norm by choosing a proper Wasserstein distance. A significant contribution of Ao et al. (2018), on the other hand, is that it allows a large size of the underlying portfolio, even although it must be of the same order of the sample data size.

The rest of the paper is organized as follows: In Section 2 we formulate the DRMV model and present necessary preliminaries. Section 3 demonstrates the tractability of our DRMV model after a series of transformations, and Section 4 studies the choices of distributional uncertainty size and the worst return level. Then, in Section 5, we report the empirical performance of our strategies against those of several other models. Concluding remarks are given in Section 6. Technical proofs of our results are given in appendices at the end of the paper.

## 2 Model Formulation

In this section, we formulate our distributionally robust Markowitz (DRM) model while reviewing some useful concepts.

Let  $\mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d)$  be the space of all Borel probability measures supported on  $\mathbb{R}^d \times \mathbb{R}^d$ . A given element  $\pi \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d)$  can be assumed to be the joint distribution of a random vector  $(U, V)$ , where  $U \in \mathbb{R}^d$  and  $V \in \mathbb{R}^d$ . We use  $\pi_U$  and  $\pi_V$  to denote the marginal distributions of  $U$  and  $V$  under  $\pi$ . In particular,  $\pi_U(A) = \pi(A \times \mathbb{R}^d)$  and  $\pi_V(A) = \pi(\mathbb{R}^d \times A)$  for every Borel set  $A \subset \mathbb{R}^d$ .

We start with a “cost” function  $c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, \infty]$ , which we shall assume to be lower semicontinuous and such that  $c(u, u) = 0$  for any  $u \in \mathbb{R}^d$ . For a given such cost function  $c$ , we

---

<sup>4</sup>Cross validation, while a standard technique, is generally data intensive and time consuming. In the rolling horizon setting, for instance, one has to assume the parameter  $\delta$  to be constant during the horizon and estimate it. In contrast, in our data-driven approach, the parameter  $\delta$  changes over the horizon in a way that is sensitive to the variability in the data.

introduce  $D_c(\cdot, \cdot)$  representing some “discrepancy” between two probability measures as follows:

$$D_c(\mathbb{P}, \mathbb{Q}) := \inf\{\mathbb{E}_\pi[c(U, V)] : \pi \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d), \pi_U = \mathbb{P}, \pi_V = \mathbb{Q}\},$$

where  $\mathbb{P}$  and  $\mathbb{Q}$  are two probability measures supported on  $\mathbb{R}^d$ . This can be interpreted as the optimal (minimal) transportation cost (also known as the optimal transport discrepancy or the Wasserstein discrepancy) of moving the mass from  $\mathbb{P}$  into the mass of  $\mathbb{Q}$  under a cost  $c(x, y)$  per unit of mass transported from  $x$  to  $y$ .

If for a given  $p > 0$ ,  $c^{1/p}(\cdot, \cdot)$  is a metric, then so is  $D_c^{1/p}(\cdot, \cdot)$ ; see Villani (2003). Such a metric  $D_c^{1/p}(\cdot, \cdot)$  is known as a Wasserstein distance of order  $p$ . Most of the times in this paper, we choose the following cost function

$$c(u, v) = \|u - v\|_q^2$$

where  $q \geq 1$  is fixed (which leads to a Wasserstein distance of order 2).<sup>5</sup>

Recall that  $R$  is the  $d$ -dimensional vector of random returns of the  $d$  stocks. Let  $\mathbb{P}_n$  be the empirical probability measure on  $\mathbb{R}^d$  with a sample size  $n$ , i.e

$$\mathbb{P}_n(dx) = \frac{1}{n} \sum_{i=1}^n \delta_{R_i}(dx)$$

where  $R_i$  ( $i = 1, 2, \dots, n$ ) are realizations of  $R$  and  $\delta_{R_i}(\cdot)$  is the indicator function. Define the ambiguity set as

$$\mathcal{U}_\delta(\mathbb{P}_n) = \{\mathbb{P} : D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta\},$$

and the feasible region of portfolios as

$$\mathcal{F}_{\delta, \bar{\alpha}}(n) = \{\phi \in \mathbb{R}^d : \phi^\top \mathbf{1} = 1, \min_{\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n)} [\mathbb{E}_\mathbb{P}(\phi^\top R)] \geq \bar{\alpha}\}.$$

The DRMV approach then consists in choosing a portfolio  $\phi \in \mathcal{F}_{\delta, \bar{\alpha}}(n)$  which achieves the optimal min-max value in (2).

### 3 Transformations, Duality and Regularization

Problem (2) appears, in principle, very complex. First of all, the inner maximization problem is over a set of probability measures, which renders it an infinite dimensional optimization problem. Second, it is not clear whether the outer minimization problem, while finite dimensional, is convex. Therefore (2) at its outset seems computational insurmountable. In this section, we reformulate (2), through a series of transformations and a duality argument, as an equivalent problem that is computationally tractable.

The following theorem, whose proof is relegated to Appendix A, is the main result of the paper, which states that (2) is equivalent to a non-robust portfolio selection problem in terms

---

<sup>5</sup>Different cost functions can be used, resulting in different regularization penalties, as we will discuss at the end of Section 3 and in Section 6.

of the empirical measure  $\mathbb{P}_n$ , with an additional “regularization” term.

**Theorem 1** *The primal formulation given in (2) is equivalent to the following dual problem*

$$\begin{cases} \min_{\phi \in \mathbb{R}^d} & \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} + \sqrt{\delta} \|\phi\|_p, \\ \text{subject to} & \phi^\top \mathbf{1} = 1, \mathbb{E}_{\mathbb{P}_n}(\phi^\top R) \geq \bar{\alpha} + \sqrt{\delta} \|\phi\|_p, \end{cases} \quad (3)$$

*in the sense that the two problems have the same optimal solutions and optimal value.*

It is not hard to verify that the mapping  $\phi \rightarrow \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} + \sqrt{\delta} \|\phi\|_p$  is convex, and the feasible region of (3) is clearly convex. So (3) and therefore (2) are both convex optimization problems. As such, they are tractable optimization problems.

Problem (3) has an additional term,  $\sqrt{\delta} \|\phi\|_p$ , in its objective function. In asset management industry, fund managers using mean–variance portfolio selection model often add a “penalty” or “regularization” term – in the form of  $k \|\phi\|_p$  where  $\|\cdot\|_p$  is an appropriately chosen norm – in order to enhance the sparsity of the vector as a way to include less stocks to the portfolios and to address the issue of overfitting.<sup>6</sup> Here, we provide *interpretability* of this regularization technique (which is based on experience or heuristics) by a well-established robustification idea, backed by precise rationality principles; see for example Delage et al. (2019). Moreover, the parameter  $\delta$  that reflects the level of regularization will also be endogenously informed by data, as we will show in the next section.

Theorem 1 has another interesting interpretation related to transactions costs. Introducing the Lagrange multiplier to the second constraint in Problem (3), the latter is equivalent to

$$\begin{cases} \min_{\phi \in \mathbb{R}^d} & \gamma \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} - \mathbb{E}_{\mathbb{P}_n}(\phi^\top R) + k \|\phi\|_p, \\ \text{subject to} & \phi^\top \mathbf{1} = 1 \end{cases} \quad (4)$$

for some constants  $\gamma > 0$  and  $k \geq 0$ . Following Olivares-Nadal and DeMiguel (2018), we can regard the term in the objective function,  $k \|\phi\|_p$ , as a transaction cost term. Therefore, (4) is a classical mean–variance model with transaction costs.<sup>7</sup> With this interpretation, Theorem 1 yields that a mean–variance model with transaction costs in the form of (4) is equivalent to a *distributionally* robust mean–variance model in the form of (2). This result is related to one of the results in Olivares-Nadal and DeMiguel (2018) which states that a mean–variance portfolio problem with  $L^p$ -norm transaction costs is equivalent to a robust optimization problem. However, there are important differences between the two results. The robust problem in Olivares-Nadal and DeMiguel (2018) has an *ellipsoidal* uncertainty set around the sample mean *only*, and there is no robustification on the variance. The distributional robust model in our paper is the most comprehensive one, because it robustifies not only the mean but also the variance, and indeed the whole distribution. Methodologically, the result of Olivares-Nadal

<sup>6</sup>In practice, it is not desirable to include, in the case of S&P 500 stocks for example, all the 500 stocks in one’s portfolios, even though one of the key implications of the mean–variance model is diversification. From a practical perspective including too many stocks is costly and prone to mismanagement. Therefore, adding a proper regularization term not only reduces overfitting, but also helps achieve a balance between diversification and manageability.

<sup>7</sup>Strictly speaking, (4) is a mean–standard-deviation model, which is equivalent to the mean–variance one.

and DeMiguel (2018) follows directly (and indeed trivially) from what is essentially the Legendre transformation or the convex conjugate (see the online companion of Olivares-Nadal and DeMiguel (2018)). This implication was actually well documented in earlier papers such as Bertsimas et al. (2004) and, in a portfolio selection context, Gotoh and Takeda (2011), which shows that any norm constraint can be turned into a robust constraint associated with the return vector. In contrast, the proof of Theorem 1 is much more involved and subtle because of the distributional constraint (see Appendix A). In turn, the distributional formulation is key in providing a natural statistical approach to selecting the size of the uncertainty. These constitute one of the main methodological contributions of this paper.

As indicated earlier, it should be noted that Theorem 1 does not directly follow from any of the strong duality results mentioned in the introduction. This is because the portfolio variance in the objective function (2) is *not* a linear function of the probability measure. A related work, Gao et al. (2017) proves only an asymptotic equivalence to regularization. On the other hand, we have the exact equivalence between (2) and a regularized optimization problem given in Theorem 1.

To conclude this section, we note that while the cost function is chosen as  $c(u, v) = \|u - v\|_q^2$  in the study here, our result (Theorem 1) actually holds for any cost function of the form  $c(u, v) = \|u - v\|^2$  where  $\|\cdot\|$  is any given norm with a suitable dual. More precisely, define the dual norm as  $\|x\|_* = \sup_{\|z\|=1} |x^\top z|$ . Then the primal distributionally robust model under this alternative cost function is equivalent to the following dual problem:

$$\min_{\phi \in \mathcal{F}_{\delta, \bar{\alpha}}(n)} \left( \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} + \sqrt{\delta} \|\phi\|_* \right)^2,$$

where the feasible region is modified as

$$\mathcal{F}_{\delta, \bar{\alpha}}(n) = \{\phi \in \mathbb{R}^d : \phi^\top \mathbf{1} = 1, \mathbb{E}_{\mathbb{P}_n}(\phi^\top R) \geq \bar{\alpha} + \sqrt{\delta} \|\phi\|_*\}.$$

For example, consider a norm as  $\|x\| = (x^\top \Sigma x)^{1/2}$  where  $\Sigma$  is a strictly positive definite matrix. Then  $\|x\|_* = (x^\top \Sigma^{-1} x)^{1/2}$ . Interested readers may refer to Blanchet and Kang (2017) for discussions on some other interesting norms.

## 4 Choice of Model Parameters

There are two key parameters,  $\delta$  and  $\bar{\alpha}$ , in the formulation (2), the choice of which is not only curious in theory, but also crucial in practical implementation and for the success of our algorithm. The idea is that the choice of these parameters should be informed by the data (i.e. in a data-driven way) based on some statistical principles, rather than being arbitrarily exogenous. Specifically, we define the distributional uncertainty region just large enough so that the correct optimal portfolio (the one which we would apply if the underlying distribution was known) becomes a plausible choice with a sufficiently high confidence level. Once this is determined, then we determine the feasible set of portfolios just large enough so that the correct optimal portfolio is feasible with adequately high confidence.



We need to impose several technical/statistical assumptions.

**A1)** The underlying return time series  $(R_k : k \geq 0)$  is a stationary, ergodic process satisfying  $\mathbb{E}_{\mathbb{P}^*} (\|R_k\|_2^4) < \infty$  for each  $k \geq 0$ . Moreover, for each measurable  $g(\cdot)$  such that  $|g(x)| \leq c(1 + \|x\|_2^2)$  for some  $c > 0$ , the limit

$$\Upsilon_g := \lim_{n \rightarrow \infty} \text{Var}_{\mathbb{P}^*} \left( n^{-1/2} \sum_{k=1}^n g(R_k) \right)$$

exists and the central limit theorem holds:

$$n^{1/2} [\mathbb{E}_{\mathbb{P}_n} (g(R)) - \mathbb{E}_{\mathbb{P}^*} (g(R))] \Rightarrow N(0, \Upsilon_g),$$

where (and henceforth) “ $\Rightarrow$ ” denotes weak convergence.

**A2)** For any matrix  $\Lambda \in \mathbb{R}^{d \times d}$  and any vector  $\zeta \in \mathbb{R}^d$  such that either  $\Lambda \neq 0$  or  $\zeta \neq 0$ ,

$$\mathbb{P}^* (\|\Lambda R + \zeta\|_2 > 0) > 0.$$

**A3)** The classical model (1) has a unique solution  $\phi^*$ . Moreover,  $\text{Var}_{\mathbb{P}^*} [R]$  is positive definite.

Assumption A1) is standard for most time series models (after removing seasonality). Assumption A2) holds assuming, for example, that  $R$  has a density. Assumption A3) is a technical assumption which can be relaxed, but then the evaluation of the optimal choice of  $\delta$  would become more cumbersome, as we shall explain.

#### 4.1 Choice of $\delta$

The choice of the uncertainty size  $\delta$  is crucial. If  $\delta$  is too large, then there is too much model ambiguity and the available data becomes less relevant. In this case, the resulting optimal portfolios will tend to be just equal allocations. If  $\delta$  is too small, then the effect of robustification will be negligible. Therefore, the choice of  $\delta$  should *not* be exogenously specified; rather it should be endogenously informed by the data.

Theorem 1 actually suggests an appropriate order of  $\delta = \delta_n$  (here  $n$  is the size of the available return time series data) in terms of  $n$ . Because the differences between the optimal standard deviation by solving (1) and that obtained by solving the empirical version of (1) are of order  $O(n^{-1/2})$ , it follows from Theorem 1 that any choice of  $\delta_n$  in the order of  $o(n^{-1})$  would be too small. Hence, an “optimal” order of  $\delta_n$  should be of order  $O(n^{-1})$ .

In order to choose an appropriate  $\delta_n$ , here we follow the idea behind the RWPI approach introduced in Blanchet et al. (2016). Intuitively,  $\delta$  should be chosen such that the set  $\mathcal{U}_\delta(\mathbb{P}_n) = \{\mathbb{P} : D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta\}$  contains all the probability measures that are plausible variations of the data represented by  $\mathbb{P}_n$ . Denote by  $\mathcal{Q}(\mathbb{P})$  the classical Markowitz portfolio selection problem

with target return  $\rho$  assuming that  $\mathbb{P}$  is the underlying model:

$$\begin{aligned} \min_{1^\top \phi = 1} \quad & \phi^\top \mathbb{E}_{\mathbb{P}}[RR^\top] \phi \\ \text{subject to} \quad & \phi^\top \mathbb{E}_{\mathbb{P}}[R] = \rho, \end{aligned} \tag{5}$$

and by  $\phi_{\mathbb{P}}$  a solution to  $\mathcal{Q}(\mathbb{P})$  and  $\Phi_{\mathbb{P}}$  the set of all such solutions. According to Assumption A3) we have  $\Phi_{\mathbb{P}^*} = \{\phi^*\}$  for some portfolio  $\phi^*$ . Therefore there exist (unique) Lagrange multipliers  $\lambda_1^*$  and  $\lambda_2^*$  such that

$$\begin{aligned} 2\mathbb{E}_{\mathbb{P}^*}(RR^\top)\phi^* - \lambda_1^*\mathbb{E}_{\mathbb{P}^*}[R] - \lambda_2^*1 &= 0, \\ (\phi^*)^\top \mathbb{E}_{\mathbb{P}^*}[R] - \rho &= 0. \end{aligned} \tag{6}$$

Now, when  $\delta$  is suitably chosen so that  $\mathcal{U}_\delta(\mathbb{P}_n)$  constitutes the models that are plausible variations of  $\mathbb{P}_n$ , any  $\phi_{\mathbb{P}}$  with  $\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n)$  is a plausible estimate of  $\phi^*$ . This intuition motivates the definition of the following set

$$\Lambda_\delta(\mathbb{P}_n) = \cup_{\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n)} \Phi_{\mathbb{P}},$$

which corresponds to all the plausible estimates of  $\phi^*$ . As a result,  $\Lambda_\delta(\mathbb{P}_n)$  is a natural confidence region for  $\phi^*$  and, therefore,  $\delta$  should be chosen as the smallest number  $\delta_n^*$  such that  $\phi^*$  belongs to this region with a given confidence level. Namely,

$$\delta_n^* = \min\{\delta > 0 : \mathbb{P}^*(\phi^* \in \Lambda_\delta(\mathbb{P}_n)) \geq 1 - \delta_0\},$$

where  $1 - \delta_0$  is a user-defined confidence level (typically 95%).

However, by the mere definition, it is difficult to compute  $\delta_n^*$ . We now provide a simpler representation for  $\delta_n^*$  via an auxiliary function called the robust Wasserstein profile (RWP) function. To this end, first observe that any  $\phi \in \Lambda_\delta(\mathbb{P}_n)$  if and only if there exist  $\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n)$  and  $\lambda_1, \lambda_2 \in (-\infty, \infty)$  such that

$$\begin{aligned} 2\mathbb{E}_{\mathbb{P}}(RR^\top)\phi - \lambda_1\mathbb{E}_{\mathbb{P}}[R] - \lambda_21 &= 0, \\ \phi^\top \mathbb{E}_{\mathbb{P}}(R) - \rho &= 0. \end{aligned}$$

From these two equations, multiplying the first equation by  $\phi$ , substituting the expression in the second equation and noting that  $\phi \cdot 1 = 1$ , we obtain

$$\lambda_2 = 2(\phi)^\top \mathbb{E}_{\mathbb{P}}(RR^\top)\phi - \lambda_1\rho.$$

We now define the following RWP function

$$\bar{\mathcal{R}}_n(\phi, \lambda_1, \Sigma, \mu) := \inf \left\{ D_c(\mathbb{P}, \mathbb{P}_n) : \left\{ \begin{array}{l} 2\Sigma\phi - \lambda_1\mu = \left( 2(\phi)^\top \Sigma\phi - \lambda_1\mu \cdot \phi \right) 1 \\ \mu = \mathbb{E}_{\mathbb{P}}[R], \Sigma = \mathbb{E}_{\mathbb{P}}(RR^\top) \end{array} \right. \right\},$$

for  $(\phi, \lambda_1, \Sigma, \mu) \in \mathbb{R}^d \times \mathbb{R} \times \mathcal{S}_+^{d \times d} \times \mathbb{R}^d$  where  $\mathcal{S}_+^{d \times d}$  is the set of all the symmetric positive semidefinite matrices and we convent that  $\inf \emptyset := +\infty$ . Moreover, define

$$\bar{\mathcal{R}}_n^*(\phi^*) := \inf_{\Sigma \in \mathcal{S}_+^{d \times d}, \mu \in \mathbb{R}^d, \lambda_1 \in \mathbb{R}} \bar{\mathcal{R}}_n(\phi^*, \lambda_1, \Sigma, \mu).$$

It follows directly from the definitions that

$$\phi^* \in \Lambda_\delta(\mathbb{P}_n) \implies \bar{\mathcal{R}}_n^*(\phi^*) \leq \delta, \quad (7)$$

while for any given  $\epsilon > 0$

$$\bar{\mathcal{R}}_n^*(\phi^*) \leq \delta + \epsilon \implies \phi^* \in \Lambda_\delta(\mathbb{P}_n). \quad (8)$$

Let us define

$$\tilde{\delta}_n^* = \inf\{\delta > 0 : \mathbb{P}^*(\bar{\mathcal{R}}_n^*(\phi^*) \leq \delta) \geq 1 - \delta_0\}.$$

It follows from (7) and (8) that  $\tilde{\delta}_n^* \leq \delta_n^* \leq \tilde{\delta}_n^* + \epsilon$ . As  $\epsilon > 0$  is arbitrary, we obtain

$$\delta_n^* = \tilde{\delta}_n^* = \inf\{\delta : \mathbb{P}^*(\bar{\mathcal{R}}_n^*(\phi^*) \leq \delta) \geq 1 - \delta_0\}.$$

In other words,  $\delta_n^*$  is the quantile corresponding to the  $1 - \delta_0$  percentile of the distribution of  $\bar{\mathcal{R}}_n^*(\phi^*)$ .<sup>8</sup>

Still, even under A3), the statistic  $\bar{\mathcal{R}}_n^*(\phi^*)$  is somewhat cumbersome to work with as it is derived from solving a minimization problem in terms of the mean and variance. So, instead, we will define an alternative statistic involving only the empirical mean and variance while producing an upper bound of  $\delta = \delta_n$  which still preserves the target rate of convergence to zero as  $n \rightarrow \infty$  (which, as we have argued, should be of order  $O(n^{-1})$ ).

Denote  $\Sigma_n = \mathbb{E}_{\mathbb{P}_n}(RR^\top)$ , and let  $\lambda_1^*$  be the Lagrange multiplier in (6). Set

$$\mu_n = \rho \mathbf{1} + 2(\Sigma_n \phi^* - \phi^{*T} \Sigma_n \phi^* \mathbf{1}) / \lambda_1^*. \quad (9)$$

Define

$$\mathcal{R}_n(\Sigma_n, \mu_n) := \bar{\mathcal{R}}_n(\phi^*, \lambda_1^*, \Sigma_n, \mu_n).$$

It is clear that

$$\mathcal{R}_n(\Sigma_n, \mu_n) \geq \bar{\mathcal{R}}_n^*(\phi^*).$$

Therefore,

$$\mathcal{R}_n(\Sigma_n, \mu_n) \leq \delta \implies \bar{\mathcal{R}}_n^*(\phi^*) \leq \delta$$

and, consequently,

$$\bar{\delta}_n^* = \inf\{\delta \geq 0 : \mathbb{P}^*(\mathcal{R}_n(\Sigma_n, \mu_n) \leq \delta) \geq 1 - \delta_0\} \geq \delta_n^*. \quad (10)$$

---

<sup>8</sup>Herein the analysis is under Assumption A3). If  $\Phi_{\mathbb{P}^*}$  contained more than just one element, then there would be several possible options to formulate an optimization problem for choosing  $\delta$ . For example, we may choose  $\delta$  as the smallest uncertainty size such that  $\Phi_{\mathbb{P}^*} \subset \Lambda_\delta(\mathbb{P}_n)$  with probability  $1 - \delta_0$ , in which case we would need to study  $\sup_{\phi^* \in \Phi_{\mathbb{P}^*}} \bar{\mathcal{R}}_n^*(\phi^*)$ .

Moreover, because of the choice of  $\Sigma_n$  and  $\mu_n$ , we have

$$\mathcal{R}_n(\Sigma_n, \mu_n) = \inf\{\mathcal{D}_c(\mathbb{P}, \mathbb{P}_n) : \mathbb{E}_{\mathbb{P}}[RR^\top] = \Sigma_n, \mathbb{E}_{\mathbb{P}}[R] = \mu_n\}.$$

The next result shows  $\bar{\delta}_n^* = O(n^{-1})$  as  $n \rightarrow \infty$ .

**Theorem 2** *Assume A1) and A2) hold and write  $\mu_* = \mathbb{E}_{\mathbb{P}^*}(R)$  and  $\Sigma_* = \mathbb{E}_{\mathbb{P}^*}(RR^\top)$ . Define  $g(x) = x + 2(xx^\top \cdot \phi^* - \phi^{*T}xx^\top\phi^*1)/\lambda_1^*$ . Then*

$$n\mathcal{R}_n(\Sigma_n, \mu_n) \Rightarrow L_0 := \sup_{\bar{\lambda} \in \mathbb{R}^d} \left( \bar{\lambda}^\top Z - \inf_{\bar{\Lambda} \in \mathbb{R}^{d \times d}} \mathbb{E}_{\mathbb{P}^*}[\|\bar{\Lambda}R + \bar{\lambda}\|_p^2] \right)$$

where  $Z \sim N(0, \Upsilon_g)$ . Moreover, if  $p = 2$  then

$$L_0 = \frac{\|Z\|_2^2}{4(1 - \mu_*^\top \Sigma_*^{-1} \mu_*)}.$$

A proof of Theorem 2 is provided in Appendix B.

Note that  $L_0$  has an explicit expression when  $p = 2$ . When  $p \neq 2$ , using the inequalities that  $\|x\|_p^2 \geq \|x\|_2^2$  if  $p < 2$  and  $d^{(\frac{1}{2}-\frac{1}{p})}\|x\|_p^2 \geq \|x\|_2^2$  if  $p > 2$ , we can find a stochastic upper bound of  $L_0$  that can be explicitly expressed. In that case we can obtain  $\bar{\delta}_n^*$  in exactly the same way, namely, first compute the  $1 - \delta_0$  quantile of  $L_0$  and then let  $\bar{\delta}_n^*$  be such quantile multiplied by  $1/n$ . The distribution of  $L_0$  can be calibrated using a natural plug-in estimator, leading to an asymptotically equivalent estimator of  $\bar{\delta}_n^*$ . The validity of this type of (plug-in) approach is explained in the next section in the context of choosing  $\bar{\alpha}$ , but the principle applies directly in the setting of  $L_0$  as well. In simple words, whenever an asymptotic limiting distribution depends continuously on various parameters and consistent estimators are available for those parameters, then consistent plug-in estimators can be safely used and still preserving exactly the same asymptotic distributions. The details of this approach are investigated in Proposition 2 of Blanchet et al. (2019), and the performance of such plug-in estimators are tested empirically in Section 5 of this paper.

## 4.2 Choice of $\bar{\alpha}$

Once  $\delta$  has been chosen, the next step is to choose  $\bar{\alpha}$ . The idea is to select  $\bar{\alpha}$  just large enough to make sure that we do not rule out the inclusion  $\phi^* \in \mathcal{F}_{\delta, \bar{\alpha}}(n)$  with a given confidence level chosen by the user, where  $\phi^*$  is the optimal solution to (1). It is equivalent to choosing  $v_0$  where

$$\bar{\alpha} = \rho - \sqrt{\delta} \|\phi^*\|_p v_0.$$

Therefore, it follows from Proposition 2 (in Appendix A) that  $\phi^* \in \mathcal{F}_{\delta, \bar{\alpha}}(n)$  if and only if

$$(\phi^*)^\top \mathbb{E}_{\mathbb{P}_n}(R) - \sqrt{\delta} \|\phi^*\|_p \geq \rho - \sqrt{\delta} \|\phi^*\|_p v_0.$$

However,  $\rho = (\phi^*)^\top \mathbb{E}_{\mathbb{P}^*}(R)$ ; so the previous inequality holds if and only if

$$(\phi^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)] \geq \|\phi^*\|_p \sqrt{\delta} (1 - v_0). \quad (11)$$

Hence, we can choose  $\sqrt{\delta} (1 - v_0) < 0$  sufficiently negative so that the previous inequality holds with a specified confidence level. We hope to choose a  $v_0$  such that  $\phi^*$  will satisfy (11) with confidence level  $1 - \epsilon$ . This can be achieved asymptotically by a central limit theorem as the following result indicates.

**Proposition 1** *Suppose that A1) and A3) hold and let  $\{\phi_n^*\}_{n=1}^\infty$  be any consistent sequence of estimators of  $\phi^*$  in the sense that  $\phi_n^* \rightarrow \phi^*$  in probability as  $n \rightarrow \infty$ . Then,*

$$n^{1/2} \left\{ \frac{(\phi_n^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi_n^*\|_p} \right\} \Rightarrow N(0, \Upsilon_{\phi^*})$$

as  $n \rightarrow \infty$ , where

$$\Upsilon_{\phi^*} := \lim_{n \rightarrow \infty} \text{Var}_{\mathbb{P}^*} \left( n^{-1/2} \sum_{k=1}^n (\phi^*)^\top R_k / \|\phi^*\|_p \right).$$

A proof of this proposition is delayed to Appendix C.

Using the previous result we can estimate  $v_0$  asymptotically. Let  $\phi_n$  denote the optimal solution of problem  $\mathcal{Q}(\mathbb{P}_n)$ . We know that  $\phi_n$  converges to  $\phi^*$  in probability. So, we choose a  $v_0$  such that the following inequality will hold with confidence level  $1 - \epsilon$ ,

$$\frac{1}{\|\phi_n\|_p} (\phi_n)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)] \geq \sqrt{\delta} (1 - v_0). \quad (12)$$

According to Proposition 1 the left-hand side of (12) is approximately normally distributed and thus we can choose its  $1 - \epsilon$  quantile and consequently decide the value of  $v_0 > 1$ .

For reader's convenience, we present a simple "menu" for estimating  $\delta$  and  $\bar{\alpha}$ .

1. Choose the target return rate  $\rho$ .
2. Collect return data  $\{R_i\}_{i=1}^n$ .
3. Use the sample mean  $\mu_n = \mathbb{E}_{\mathbb{P}_n}(R)$  and the sample second-moment matrix  $\Sigma_n = \mathbb{E}_{\mathbb{P}_n}(RR^\top)$  to approximate  $\mu_*$  and  $\Sigma_*$ , respectively, appearing in Theorem 2.
4. Use the solution  $\phi_n$ , which is the solution to problem  $\mathcal{Q}(\mathbb{P}_n)$  (see (5)), to approximate  $\phi^*$  in Theorem 2.
5. Apply Theorem 2 and (10) to determine  $\delta = \bar{\delta}^*$  with the 95% confidence level.
6. Choose  $v'_0$  based on the 95% quantile according to (12) and Proposition 1, and consequently obtain  $\bar{\alpha}$ . Choose  $v''_0$  such that the equation  $(\phi_n)^\top \mathbb{E}_{\mathbb{P}_n}(R) - \sqrt{\delta} \|\phi_n\|_p = \rho - \sqrt{\delta} \|\phi_n\|_p v''_0$  holds. Then set  $v_0 = \max(v'_0, v''_0)$ .

To conclude this section, we note that the choices of  $\delta$  and  $\alpha$  are separate, each with a given confidence level. Jointly, the chosen  $(\delta, \alpha)$  may have a completely different confidence level. It remains an interesting problem to develop a joint way of choosing the two parameters to ensure a *given, fixed* confidence level.

## 5 Empirical Performance and Comparisons

In this section we report the results of our backtesting experiments on S&P 500 constituents that compare the performance of our DRMV portfolios with those of the portfolios based on the following models: classical (non-robust) single-period Markowitz, continuous-time Markowitz, Fama–French, Black–Litterman, robust Goldfarb–Iyengar, Olivares–Nadal–DeMiguel, and an equally weighted portfolio. The first four models are well-known and have been widely used in practice, the fifth one is an alternative robust model not based on distributional uncertainty, and the sixth one has transaction costs turned into an equivalent robust model that has an ellipsoidal uncertainty around the mean. The equally weighted strategy is actually an extreme outcome of the DRMV model when the uncertainty size  $\delta = \infty$ .

### 5.1 Experiment design and data preparation

We backtested for the period January 2000–December 2016 with the training (estimation) period being January 1991–December 1999 (i.e. the previous 10 years).<sup>9</sup> All the stock monthly price data had been obtained from the database of Wharton Business School. At the beginning of 2000, we *randomly* chose 100 stocks from the constituents of S&P 500 that have at least 10 years’ historical price data available.<sup>10</sup> The basic period is set to be one year in all the single-period models involved with target annual mean return rate fixed to be  $\rho = 10\%$  where applicable. Then we used the training period to estimate the out-of-sample parameters, namely the mean and the variance, to construct the optimal strategies of the various models tested.

#### 5.1.1 DRMV model

Let us first describe in details the construction of the DRMV strategy for the selected 100 stocks. We generated this 17-year long strategy in an (overlapping) rolling horizon fashion, with each horizon being a month. Specifically, on the first trading day of January 2000, we solved our DRMV model to obtain a portfolio, denoted as  $\phi_R$ . In doing so, we set  $p = q = 2$  and

---

<sup>9</sup>We chose the period 2000–2016 for our backtestings for a reason: the market was overall very volatile during this period, experiencing two major crashes: the dot com bubble burst and the subprime financial crisis, followed by a long bull run until this day of writing (February 2018). We were particularly interested to see how “robust” our DRMV strategies would have been when sailing through such a bumpy journey.

<sup>10</sup>In theory, we should have included *all* the constituents of S&P 500 in our portfolios. However, that would be computationally inefficient and practically (almost) infeasible for most of the models under testing (e.g. the original Markowitz model). Therefore, it is desirable to choose a small subset of stocks based on which to apply various models. This “stock selection” is an ultimately important part of the overall portfolio management. In this paper, however, we aim to test the performances of “stock allocation” (namely, to allocate wealth among the stocks that have been *already* selected in order to achieve the best risk-adjusted return) of these models. That is why we randomly selected the small subset of stocks in order to focus on the part of the *stock allocation*. On the other hand, the requirement that the selected stocks have at least 10 years’ price data is due to the length of the training period.

$\rho = 10\%$ , and obtained the parameters,  $\delta$  and  $\bar{\alpha}$ , using the menu at the end of Section 4. We then substituted  $\delta$  and  $\bar{\alpha}$  in the optimization problem described in Theorem 1 to obtain  $\phi_R$ .

We kept  $\phi_R$  until only the first trading day of February 2000. At that point we re-estimated the parameters  $\delta$  and  $\bar{\alpha}$  using the *immediate* previous 10-year (namely February 1991 – January 2000) price data and re-solved the DRMV model, and generated a new portfolio  $\phi_R$  for February 2000, the second month in our backtesting period. We repeated the same steps for all the subsequent months.

If at the beginning of a month, some stocks in our portfolio had been removed out of the S&P 500 during the previous month, then we would also remove them from our portfolio, replace them by the same number of stocks that were randomly picked from S&P 500 (yet having at least 10 years' historical data), and then re-balance based on our DRMV model. We still denoted by  $\phi_R$  the overall portfolio for the 17-year period and kept track of the wealth process that had been updated at the end of each month.

In what follows we will describe the implementations of the other models, mentioned at the beginning of this section, under comparison. Except for the continuous-time Markowitz model, all the rest are single-period models so we applied the same monthly rolling horizon approach to build the respective strategies. Moreover, for these single-period models, whenever there were stocks dropped from S&P 500, we would replace them with exactly the same set of stocks as in the DRMV model so as to maintain consistency across different models. The case of the continuous-time model is slightly more complicated, and we will explain how we deal with these issues separately.

### 5.1.2 Single-period and continuous-time Markowitz models

For the single-period Markowitz model, for each period (month) we used the sample mean and sample covariance matrix of the immediate previous 10-year return data to estimate the corresponding parameters in problem (1). Then we generated the optimal portfolio of the classical Markowitz model,  $\phi_M$ , by setting  $\rho = 10\%$  and solving problem (1), on exactly the same rolling horizon basis as the DRMV model.

The continuous-time Markowitz mean-variance model is based on Cui et al. (2012) in which portfolios are constructed on risky stocks only. This setting is consistent with ours.<sup>11</sup> It is assumed that the stock price process follows correlated time-inhomogeneous Black-Scholes dynamics. Let  $\{X(t) : t \in [0, T]\}$  be the wealth process (also called an admissible wealth process) under any given admissible portfolio. The mean-variance problem is

$$\text{Minimize } \text{Var}(X(T)) \tag{13}$$

subject to

$$\{X(t) : t \in [0, T]\} \text{ is admissible, } X(0) = x_0, \mathbb{E}[X(T)] = z,$$

where  $z$  is a given parameter representing the expected payoff at the end of the investment

---

<sup>11</sup>There is an extensive literature on continuous-time Markowitz models; however to our best knowledge all the other existing models include a risk-free asset.

horizon,  $T$ . An optimal strategy is given explicitly in Theorem 1.1 of Cui et al. (2012), which gives the portfolio at each given time  $t \in [0, T]$  as a function of the wealth, a couple of auxiliary feedback processes and the estimates (at time  $t$ ) of the (time-inhomogeneous) diffusion and drift coefficients.

In theory, a continuous-time model requires *continuous* rebalancing all the times. Naturally, it is not possible (indeed not necessary) in practice nor in our empirical implementation. In our experiments, we set  $T = 1$  (year) and  $z = 1.1x_0$  corresponding to an annual expected return  $\rho = 10\%$ , and we rebalance only monthly (instead of continuously). The one-year period is consistent with the other models under comparison. Therefore, on the first trading day of January 2000, we estimated all the necessary parameters/coefficients based on the previous 10-year data and then applied the explicit formula for optimal portfolio, denoted as  $\phi_C$ , given by Theorem 1.1 of Cui et al. (2012). On the first trading day of February 2000, we applied the same formula but with updated estimates of the parameters/coefficients based on the immediate previous 10-year data. This way we have constructed a strategy for the whole year of 2000. For all the subsequent years, we repeat the same procedure to generate a 17-year long strategy  $\phi_C$  and the corresponding wealth process.

It is important to note that this model does not have an explicit no-bankruptcy constraint (i.e. it does not rule out the possibility that the wealth process may go negative during  $[0, T]$ ). Indeed, as we will see in the discussions below, this model had led to bankruptcy in *all* of our numerical experiments for portfolios of 100 stocks.<sup>12</sup> Once a bankruptcy happened, we then considered it game over and kept the zero wealth until December 2016.<sup>13</sup>

### 5.1.3 Fama-French model

The celebrated Fama-French model helps estimate the covariance matrix when the number of stocks  $d$  is close to or greater than the sample size  $n$ . In this section, we follow the approach proposed by Fan et al. (2011) to implement the Fama-French model. Specifically, we assume the stock returns follow the factor model

$$\mathbf{r} = \mathbf{B}\mathbf{f} + \mathbf{u}, \quad (14)$$

where  $\mathbf{r} = (R_1, \dots, R_d)^\top$  is the random vector of stock returns,  $\mathbf{f} = (f_1, f_2, f_3)^\top$  consists of the three factors of the Fama-French model (i.e Rm-Rf, SMB and HML) respectively,  $\mathbf{B} = (b_1, \dots, b_d)^\top$  is the vector of factor loading, and  $\mathbf{u} = (u_1, \dots, u_d)^\top$  is the vector of uncorrelated errors.

Let  $(\mathbf{f}_1, \mathbf{r}_1), \dots, (\mathbf{f}_n, \mathbf{r}_n)$  be  $n$  independent and identically distributed (i.i.d.) samples of  $(\mathbf{f}, \mathbf{r})$ .

<sup>12</sup>We have also tested for portfolios with 20 stocks and observed bankruptcy in more than half of our experiments. On the other hand, although the other six single-period models have no explicit no-bankruptcy constraint either, a total of only two instances of bankruptcy occurred in our experiments.

<sup>13</sup>Bielecki et al. (2005) solved a continuous-time mean-variance model with the no-bankruptcy constraint. However, there is a risk-free asset in that model. To have a fair comparison with the other models in which there is no risk-free account available, it is proper to choose the model of Cui et al. (2012) in our experiments. We are not aware of a work on continuous-time Markowitz model without risk-free asset and with bankruptcy prohibition.



For notational simplicity, we define

$$X = (\mathbf{f}_1, \dots, \mathbf{f}_n), \mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_n) \text{ and } \mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n),$$

where  $\mathbf{u}_i, i = 1, \dots, n$ , are the corresponding error vectors. Under model (14), the covariance matrix  $\Sigma$  of stock returns satisfies

$$\Sigma = \mathbf{B}\text{cov}(\mathbf{f})\mathbf{B}^\top + \Sigma_{\mathbf{u}}, \quad (15)$$

where  $\Sigma_{\mathbf{u}}$  is the covariance matrix of  $\mathbf{u}$ .

Then we estimate  $\Sigma$  using a substitution estimator

$$\hat{\Sigma} = \hat{\mathbf{B}}\hat{\text{cov}}(\mathbf{f})\hat{\mathbf{B}}^\top + \hat{\Sigma}_{\mathbf{u}}^\tau, \quad (16)$$

where  $\hat{\mathbf{B}} = YX^\top(XX^\top)^{-1}$ ,  $\hat{\text{cov}}(\mathbf{f}) = \frac{1}{n-1}XX^\top - \frac{1}{n(n-1)}X11^\top X^\top$ , and  $\hat{\Sigma}_{\mathbf{u}}^\tau$  is obtained by applying the same adaptive thresholding approach in Fan et al. (2011) on  $\hat{\Sigma}_{\mathbf{u}}$  with  $\hat{\Sigma}_{\mathbf{u}} = \hat{\text{cov}}(\hat{\mathbf{U}})$  being the sample covariance matrix of  $\hat{\mathbf{U}}$  with  $\hat{\mathbf{U}} = \mathbf{Y} - \hat{\mathbf{B}}\mathbf{X}$ .  $\hat{\Sigma}_{\mathbf{u}}^\tau$  is obtained by applying the same adaptive thresholding approach mentioned in Fan et al. (2011) on  $\hat{\Sigma}_{\mathbf{u}}$ . In that paper, they used the threshold  $\omega = C * 3 * \sqrt{\frac{\log d}{n}}$  where  $C = 0.1$ . However, in our experiment, the value  $C = 0.1$  leads to bankruptcy in all the cases. So we tune this parameter and find the optimal  $C = 0.01$ .

Finally, we have a substitution estimator

$$\hat{\mu} = \hat{\mathbf{B}}\bar{\mathbf{f}} \quad (17)$$

of the mean vector  $\mu$ , where  $\bar{\mathbf{f}} = \frac{1}{n} \sum_{i=1}^n \mathbf{f}_i$ .

In implementing the Fama–French model, we first downloaded the monthly data of the three factors<sup>14</sup> from Kenneth French’s personal website.<sup>15</sup> Then on the first trading day of each month during January 2000–December 2016, we used its immediate prior 10-year history returns and factors data to estimate the covariance matrix and mean vector using (16) and (17) respectively. We used these estimates for all the randomly chosen 100 stocks as the mean vector and solve the single-period classical Markowitz model with  $\rho = 10\%$ . The generated portfolio was denoted as  $\phi_F$ . This process was then repeated in the subsequent months on a rolling horizon basis.

#### 5.1.4 Black–Litterman model

The Black–Litterman model has been developed to address the mean-blur problem, namely, the fact that compared with variance, it is much more difficult to estimate within a workable accuracy the expected returns of stocks purely based on sample means. The model estimates the stock returns by the market portfolio while keeping the sample covariance matrix, and feed

<sup>14</sup>We assume that the factors have been processed according to the available papers (Fama and French (1992) and Fama and French (1993)).

<sup>15</sup> [http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data\\_library.html](http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html)

them into the classical Markowitz model to obtain the optimal strategies.

To implement this model, on the first trading day of each month during January 2000–December 2016, we calculated the implied returns of all the S&P 500 constituent stocks having at least 10 years’ historical data, by using the following formula:

$$R_{implied} = \lambda \Sigma \phi_{market},$$

where  $\lambda = 3.07$ ,  $\Sigma$  was the sample covariance matrix of the previous 10 years’ returns of these stocks and  $\phi_{market}$  was the corresponding market portfolio (i.e.,  $\phi_{market}$  is a vector whose components add up to 1, and are proportional to the capitalizations of the S&P 500 constituents having at least 10 years’ historical data) at the closing prices of the previous trading day; see Idzorek (2002).<sup>16</sup> Then we picked from  $R_{implied}$  the implied returns of the 100 stocks that had been randomly chosen. We input these returns and the sample covariance matrix into the classical Markowitz model with  $\rho = 10\%$  to obtain the portfolio  $\phi_B$ . This process was repeated in subsequent months on a rolling horizon fashion.

### 5.1.5 Goldfarb–Iyengar robust model

Goldfarb and Iyengar (2003) consider the following robust Markowitz problem with a factor model for the return rate:

$$\begin{aligned} & \text{Minimize } \phi && \max_{\mathbf{V} \in S_v, \mathbf{D} \in S_d} \text{Var}[\mathbf{r}^\top \phi] \\ & \text{subject to} && \min_{\mu \in S_m} \mathbb{E}[\mathbf{r}^\top \phi] \geq \rho, \quad \mathbf{1}^\top \phi = 1, \end{aligned}$$

where

$$\mathbf{r} = \mu + \mathbf{V}^\top \mathbf{f} + \epsilon, \quad \epsilon \sim N(0, \mathbf{D}),$$

$$S_v = \{V : V = V_0 + W, \|W_i\|_g \leq \rho_i, i = 1, \dots, n\}$$

with  $W_i$  being the  $i$ th column of  $W$  and  $\|w\|_g = \sqrt{w^\top G w}$  for some positive definite matrix  $G$ ,

$$S_d = \{D : D = \text{diag}(d), d_i \in [d_i^{min}, d_i^{max}], i = 1, \dots, n\},$$

$$S_v = \{V : V = V_0 + W, \|W_i\|_g \leq \lambda_i, i = 1, \dots, n\},$$

and

$$S_m = \{\mu : \mu = \mu_0 + \xi, |\xi_i| \leq \gamma_i, i = 1, \dots, n\}.$$

So the uncertainty set of this model is based on vector/matrix distance, as opposed to our uncertainty set that is defined through the Wasserstein distance between probability measures.

In implementing this model, we followed the instructions in Section 7.2 of Goldfarb and Iyengar (2003). Specifically, we calculated the 10 years’ sample returns  $\mathbf{r}$  of the chosen 100 stocks. Then we used the top 5 principal components of  $\mathbf{r}$  together with the return data of DJA, NDX, SPC, RUT and TYX to be the factor vector  $\mathbf{f}$ . By choosing the confidence

---

<sup>16</sup>Here we include only those having at least 10 years’ historical data to be consistent with the other models.

threshold  $\omega$  to be 95%, we estimated  $\mu_0$ ,  $\mathbf{V}_0$ ,  $\sigma_i^2$ ,  $\gamma_i$ ,  $G$  and  $\lambda_i$ . With the target annual return  $\rho = 10\%$  and plugging all the parameters from the above steps, we used SeDuMi to solve the SOCP formulation (problem (32) in Goldfarb and Iyengar (2003)) to obtain the portfolio  $\phi_G$  for each month on a rolling horizon basis, starting from January 2000.

### 5.1.6 Olivares-Nadal–DeMiguel model

Olivares-Nadal and DeMiguel (2018) examine, in the context of a mean–variance model, the equivalence between certain transaction costs and ellipsoidal robustification around the mean. They then devise a data-driven approach to portfolio selection by treating the transaction costs as a regularization term to be calibrated. Their numerical experiments test for different variants of their model, but the minimum-variance portfolios (MVPs) based on quadratic transaction costs have the overall best performance in terms of the Sharpe ratio (see Table 1 of Olivares-Nadal and DeMiguel (2018)), with which we will compare in our setting.

The MVP model is

$$\begin{aligned} & \text{Minimize } \phi && \phi^\top \Sigma \phi \\ & \text{subject to} && \mathbf{1}^\top \phi = 1 \end{aligned} \tag{18}$$

where  $\phi \in \mathbb{R}^d$  is the portfolio weight vector and  $\Sigma \in \mathbb{R}^{d \times d}$  the estimated covariance matrix of asset returns. Adding a quadratic transaction cost, Olivares-Nadal and DeMiguel (2018) consider the following model:

$$\begin{aligned} & \text{Minimize } \phi && \phi^\top \Sigma \phi + \kappa \|\Sigma^{\frac{1}{2}}(\phi - \phi_0)\|_2^2 \\ & \text{subject to} && \mathbf{1}^\top \phi = 1 \end{aligned} \tag{19}$$

where  $\kappa \in \mathbb{R}$  is the transaction cost parameter and  $\phi_0 \in \mathbb{R}^d$  is the starting portfolio. By the conjugate representation of the  $L^2$ -norm, the connection of (19) to a robust model is straightforward.

It can be shown that the optimal solution  $\phi_{ODM}$  of (19) is

$$\phi_{ODM} = \frac{1}{1 + \kappa} \phi_M + \frac{\kappa}{1 + \kappa} \phi_0, \tag{20}$$

where  $\phi_M$  solves (18).

In model (19), one needs to calibrate the parameter  $\kappa$ . Olivares-Nadal and DeMiguel (2018) do this by calibrating the trading volume  $\tau$  (i.e.,  $\|\phi - \phi_0\|_1 \leq \tau$ ). They use 10-fold cross-validation to select the best  $\tau$ . Specifically, they divide the empirical returns into 10 intervals. For each  $j$  from 1 to 10, they remove the  $j$ th interval and use the remaining returns to estimate parameters and obtain the corresponding portfolio. Then they evaluate the portfolio on the  $j$ th interval. After completing this process for each of the 10 intervals, they compute the variance of the out-of-sample returns for different  $\tau$  from the set  $\{0\%, 0.5\%, 1\%, 2.5\%, 5\%, 10\%\}$  and then choose the  $\tau$  that corresponds to the portfolio with the smallest variance.

Now we show how to infer the value of  $\kappa$  from that of  $\tau$ , once the latter has been calibrated.

By (20), the trading volume can be expressed as

$$\|\phi_{ODM} - \phi_0\|_1 = \frac{1}{1 + \kappa} \|\phi_M - \phi_0\|_1.$$

To determine  $\kappa$  we equate the trading volume to  $\tau$ , i.e.,

$$\|\phi_{ODM} - \phi_0\|_1 = \frac{1}{1 + \kappa} \|\phi_M - \phi_0\|_1 = \tau.$$

This leads to

$$\kappa = \left( \frac{\|\phi_M - \phi_0\|_1}{\tau} - 1 \right)_+,$$

where  $(x)_+ := \max(x, 0)$ .<sup>17</sup>

In implementing this model, we followed the instructions in Section 3.1 of Olivares-Nadal and DeMiguel (2018). Specifically, we calculated the 10 years' sample monthly return  $\mathbf{r}$  of the chosen 100 stocks. In the first month (January 2000), since we did not have the value of  $\phi_0$ , we set  $\kappa = 0$  and generated the optimal portfolio  $\phi_{ODM}$  of (19) for that month. In any subsequent month, we took the portfolio of the previous month as  $\phi_0$ , applied the aforementioned cross-validation method on the 10 years' sample return  $\mathbf{r}$  to select the best  $\tau$ , and plugged the corresponding  $\kappa$  into (19) to generate the optimal portfolio  $\phi_{ODM}$ .

## 5.2 Comparisons

Assume that the initial wealth at the start of the backtesting period (i.e. January 2000) is 1. For each randomly selected set of 100 stocks, we generate the wealth process for the period 2000-2016 under each of the seven models as described in the previous subsection as well as that under the equal weighting. Then we repeat the experiments on 100 such sets of 100 stocks and obtain the *average* realized wealth process for each model. These processes, along with that of S&P500 (normalized to start from 1 at the start of the testing period), are plotted in Figure 1.

Graphically, Figure 1 is “corrupted” because of the extreme behavior of the continuous-time Markowitz model. Its average performance went through the roof initially and then quickly dived to zero (all the 100 experiments ended up bankruptcy). So the continuous-time Markowitz is an extremely volatile model. This may be explained as follows. The dynamic strategies incorporate considerable feedback effects which are computed assuming the underlying model is correct. As such, model misspecifications are *compounded* precisely due to feedback effects. The inclusion of feedback in the optimal dynamic policy, in outputs which are close to typical realizations of the underlying assumed model, result in highly profitable portfolios. On the other hand, even moderate discrepancies from the underlying model dynamics might lead to relatively poor performance. As a consequence, the dynamic model exhibits significantly high variability than the static-rolling-horizon robust counterpart.

In order to be able to visualize the comparison among other portfolios, it is necessary to remove the continuous-time Markowitz from Figure 1, resulting in Figure 2. It is evident that

<sup>17</sup>Here, if  $\|\phi_M - \phi_0\|_1 \leq \tau$ , then  $\phi_M$  satisfies the trading volume constraint and hence itself is the optimal portfolio. In this case  $\kappa = 0$ .

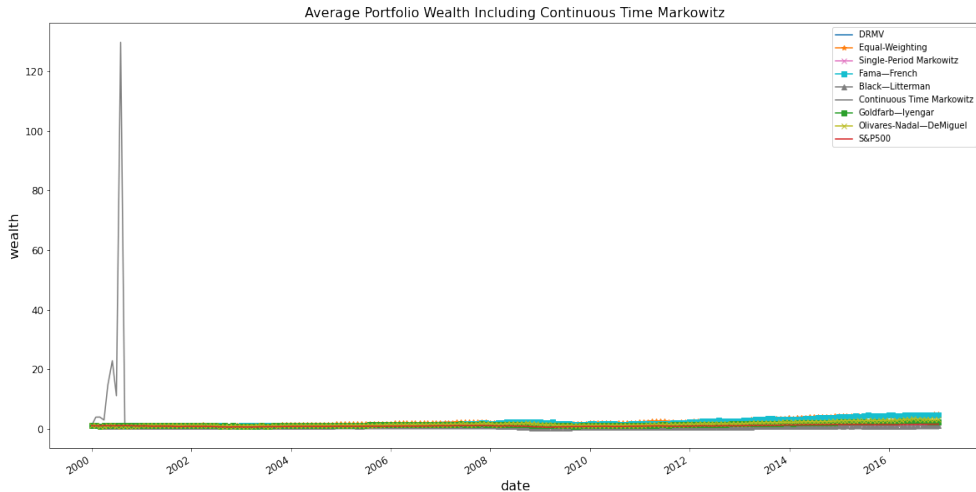


Figure 1: This graph presents wealth processes of all portfolios (including continuous time Markowitz) and S&P 500 from January 2000 to December 2016. All the portfolios except S&P 500 consist of 100 stocks and the averages are calculated over 100 numerical experiments. The  $x$ -axis indicates the time in months (from 1 to 204) and the  $y$ -axis indicates the portfolio wealth. Initial wealth is set to be 1.

all the seven models except Black–Litterman (almost) uniformly and substantially outperform S&P500 during the 17-year period. In terms of the final realized wealth, of the seven models, DRMV, equal-weighting and Fama–French lead by a substantial margin. The second-tier league includes Olivares-Nadal–DeMiguel and Goldfarb–Iyengar. The classical single-period Markowitz lags behind but still manages to outperforms the market most of the time.

The average performances of DRMV and equal-weighting are close, although the former beats the latter most of the time. This is no surprise as the latter can be regarded as an extreme case of the distributionally robust model when the uncertainty size  $\delta = \infty$ , whereas the former has a nearly “optimal” choice of  $\delta$  informed by the data.<sup>18</sup> We can study more closely the variability of the performances and the overall return–risk efficiency of the two models, by examining their histograms of annualized returns (i.e. the distributions of the annualized returns of the 100 experiments) and those of Sharpe ratios. These are plotted in Figures 3 and 4 respectively. In both figures, DRMV is more shifted to the right than equal-weighting, indicating the former outperforms the latter in the two criteria. Moreover, the two are almost equally concentrated, suggesting that both strategies have stable performance. We can also compare the histograms of kurtosis of the two; see Figure 5. There is no statistically significant differences between the two: most return distributions under both strategies are platykurtic (i.e. the kurtosis values are less than 3), implying there are fewer extreme outliers than the standard normal. Overall, we can conclude that both DRMV and equal-weighting are robust and stable, but the former is superior to the latter in terms of the return–risk efficiency.

<sup>18</sup>The good performance of the equally weighted portfolio is well documented in the literature; see e.g. DeMiguel et al. (2009).

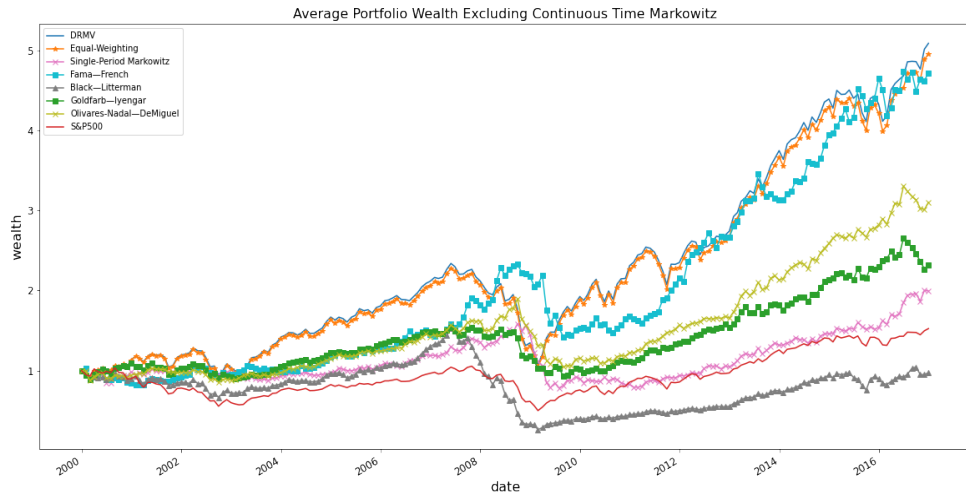


Figure 2: This graph presents wealth processes of all portfolios (excluding continuous-time Markowitz) and S&P 500 from January 2000 to December 2016. All the portfolios except S&P 500 consist of 100 stocks and the averages are calculated over 100 numerical experiments. The  $x$ -axis indicates the time in months (from 1 to 204) and the  $y$ -axis indicates the portfolio wealth. Initial wealth is set to be 1.

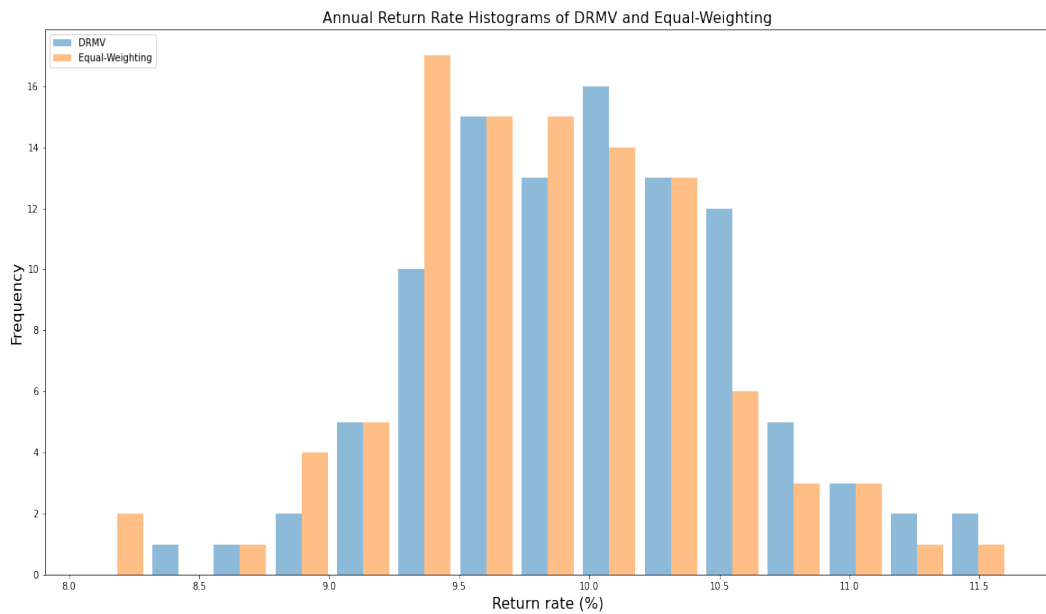


Figure 3: This graph presents the histograms of the annualized returns of the 100 different experiments on DRMV (blue) and equal-weighting (orange) portfolios. The  $x$ -axis represents the annualized returns and the  $y$ -axis represents the numbers of returns.

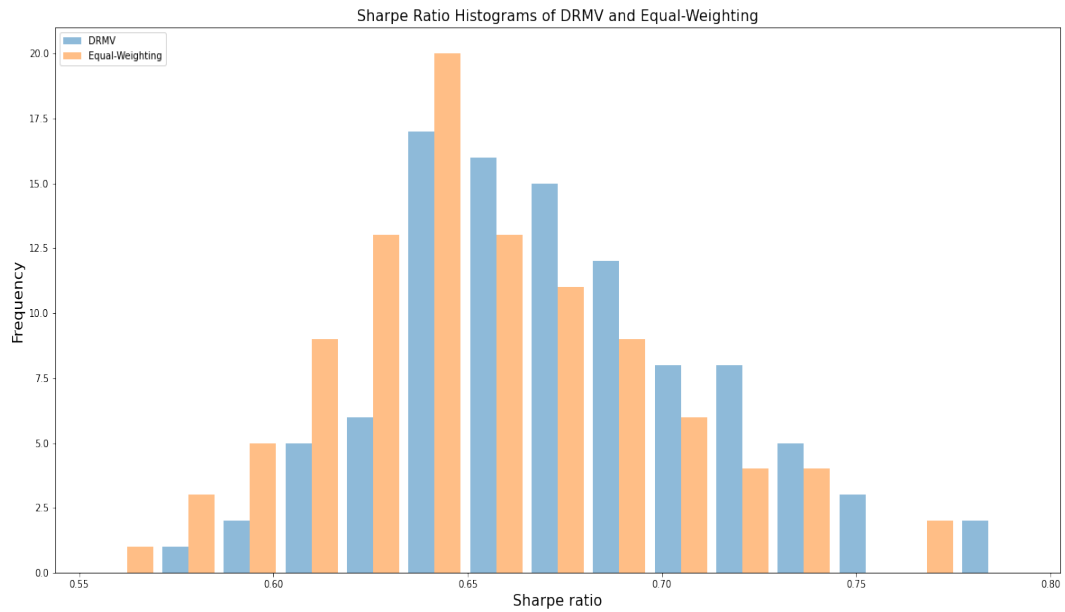


Figure 4: This graph presents the histograms of the Sharpe ratio of the 100 different experiments on DRMV (blue) and equal-weighting (orange) portfolios. The  $x$ -axis represents the Sharpe ratio and the  $y$ -axis represents the numbers of Sharpe ratios.

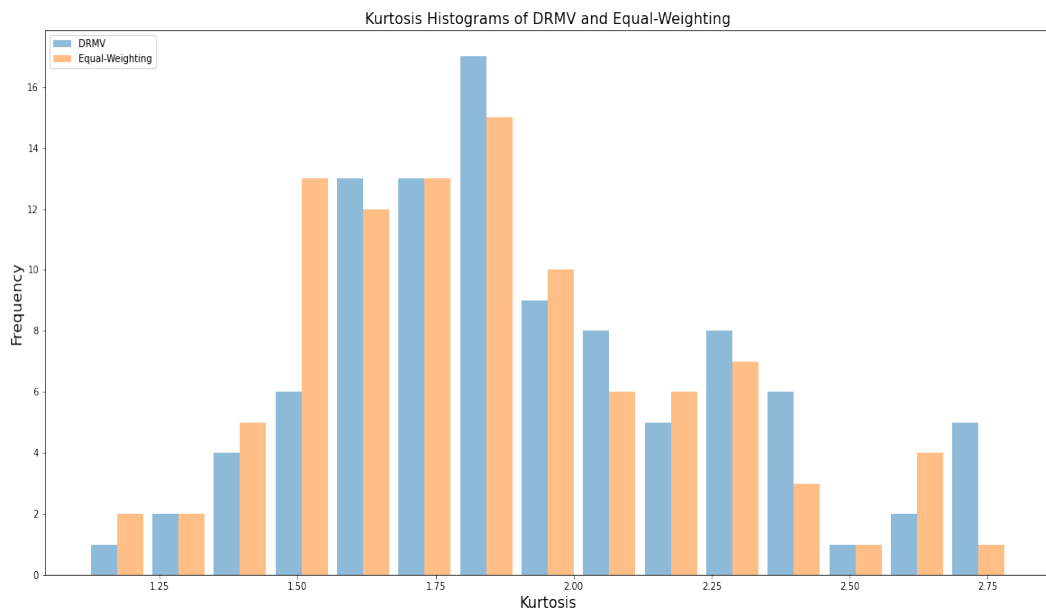


Figure 5: This graph presents the histograms of the kurtosises of the 100 different experiments on DRMV (blue) and equal-weighting (orange) portfolios. The  $x$ -axis represents the kurtosis and the  $y$ -axis represents the numbers of kurtosises.

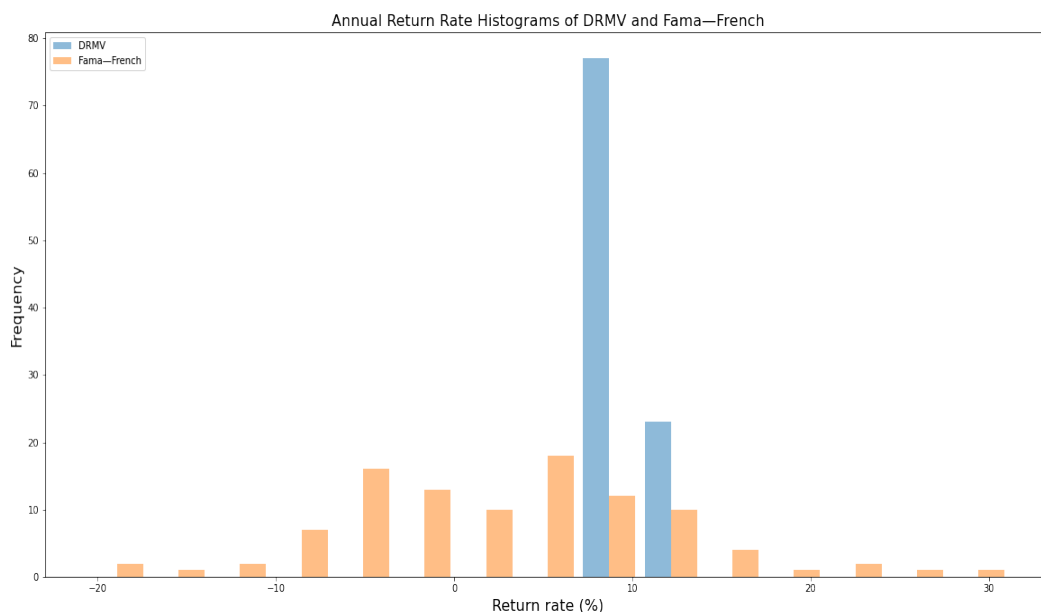


Figure 6: This graph presents the histograms of the annualized returns of the 100 different experiments on DRMV (blue) and Fama–French (orange) portfolios. The  $x$ -axis represents the annualized returns and the  $y$ -axis represents the numbers of returns.

Similarly, we compare the respective histograms between DRMV and Fama–French; see Figures 6– 8. DRMV has a much more concentrated return histogram indicating a significantly more robust performance, a much more right-shifted Sharpe ratio histogram, and a more left-shifted kurtosis histogram implying fewer extreme returns. We can therefore conclude that DRMV compares favorably with Fama–French in all the key metrics reported. However, it is interesting to note that DRMV utilizes *only* the price data, whereas Fama–French requires additional fundamental information on the companies concerned.

Likewise, DRMV outperforms the Olivares-Nadal–DeMiguel model (Olivares-Nadal and DeMiguel (2018)) based on all the histogram comparisons; see Figures 9–11. This suggests that, among other things, the inference method for selecting the uncertainty/regularization size is advantageous compared with the cross validation.

As for the robust portfolio model by Goldfarb and Iyengar (2003), we notice that it has a reasonably concentrated return histogram (Figure 12), indicating its robustness. However, DRMV’s returns are not only more concentrated, but also distributed more to the right than Goldfarb–Iyengar’s. Together with the Sharpe ratio histogram (Figure 13), the kurtosis histogram (Figure 14), and the average wealth comparison (Figure 2), it is clear that our uncertainty set formulation based on Wasserstein distance is a significant improvement to the matrix/vector distance.

Finally, we provide comparisons of histograms between DRMV and single-period Markowitz, and between DRMV and Black–Litterman; see Figures 15–20. Clearly, DRMV has far superior performance in all the metrics.



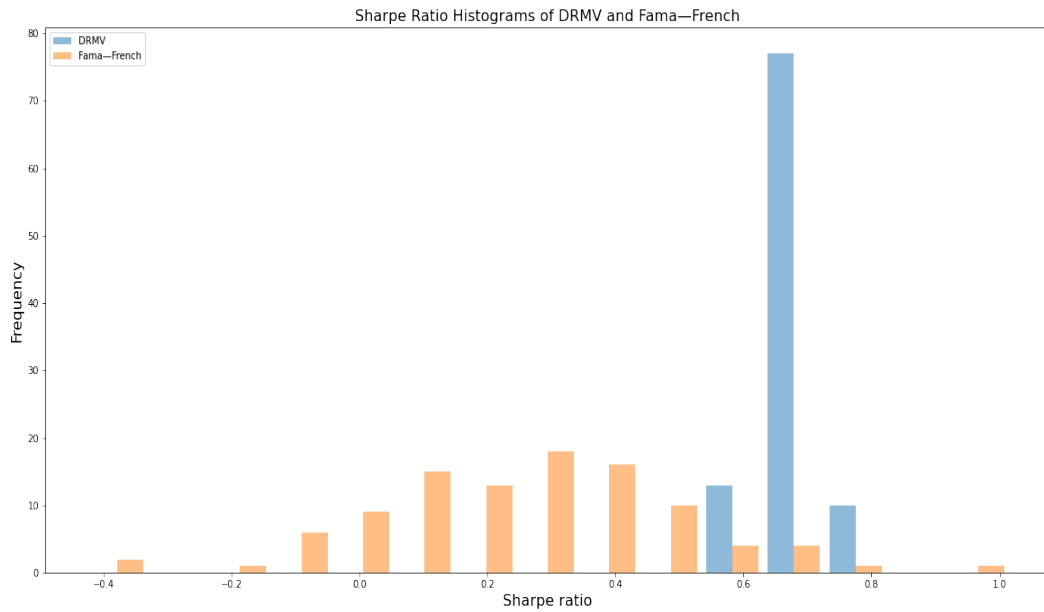


Figure 7: This graph presents the histograms of the Sharpe ratio of the 100 different experiments on DRMV (blue) and Fama-French (orange) portfolios. The  $x$ -axis represents the Sharpe ratio and the  $y$ -axis represents the numbers of Sharpe ratios.

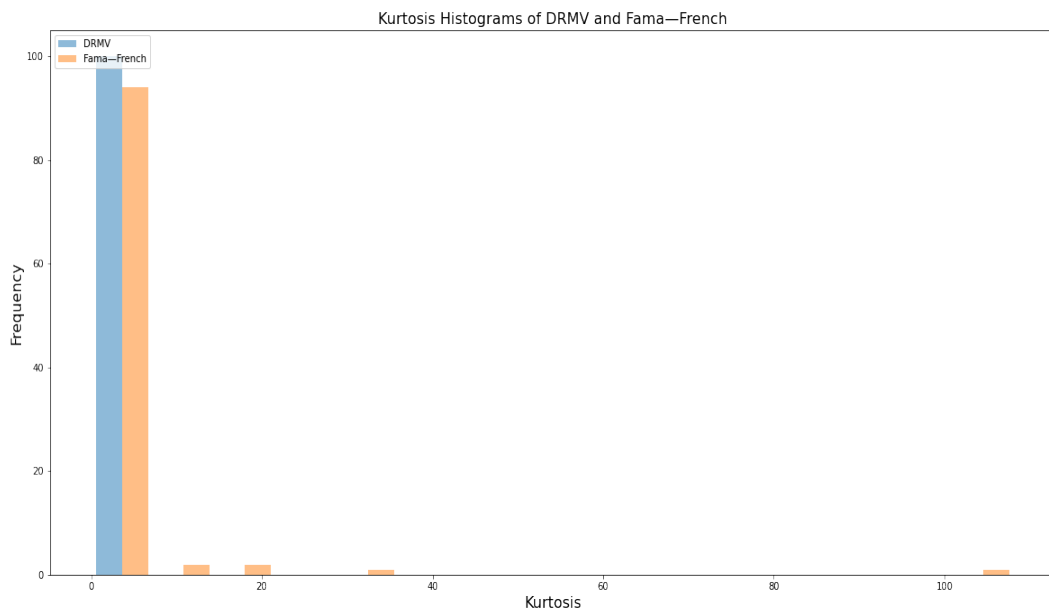


Figure 8: This graph presents the histograms of the kurtosises of the 100 different experiments on DRMV (blue) and Fama-French (orange) portfolios. The  $x$ -axis represents the kurtosis and the  $y$ -axis represents the numbers of kurtosises.

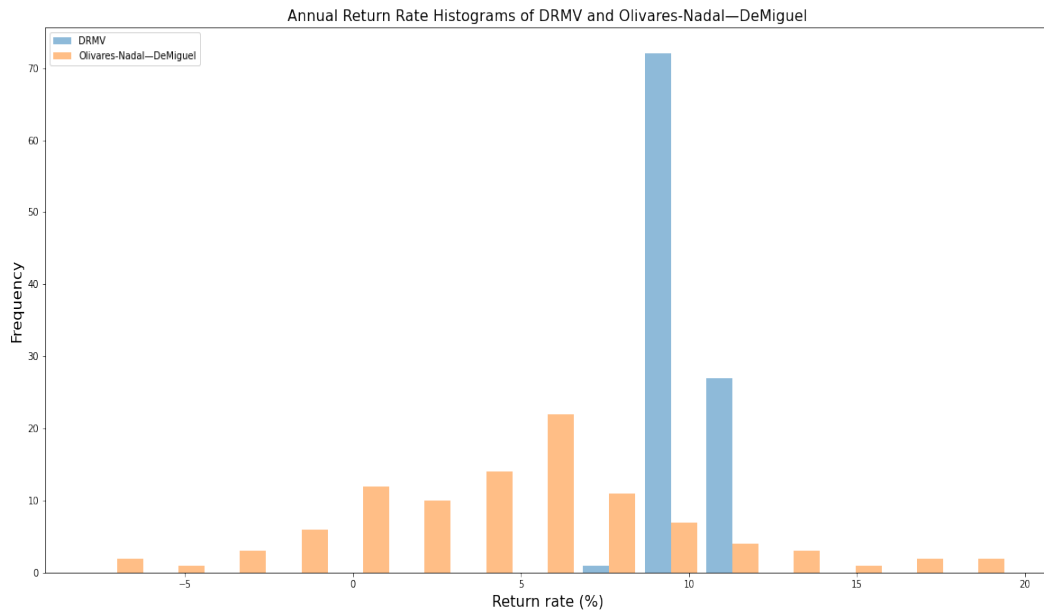


Figure 9: This graph presents the histograms of the annualized returns of the 100 different experiments on DRMV (blue) and Olivares-Nadal-DeMiguel (orange) portfolios. The  $x$ -axis represents the annualized returns and the  $y$ -axis represents the numbers of returns.

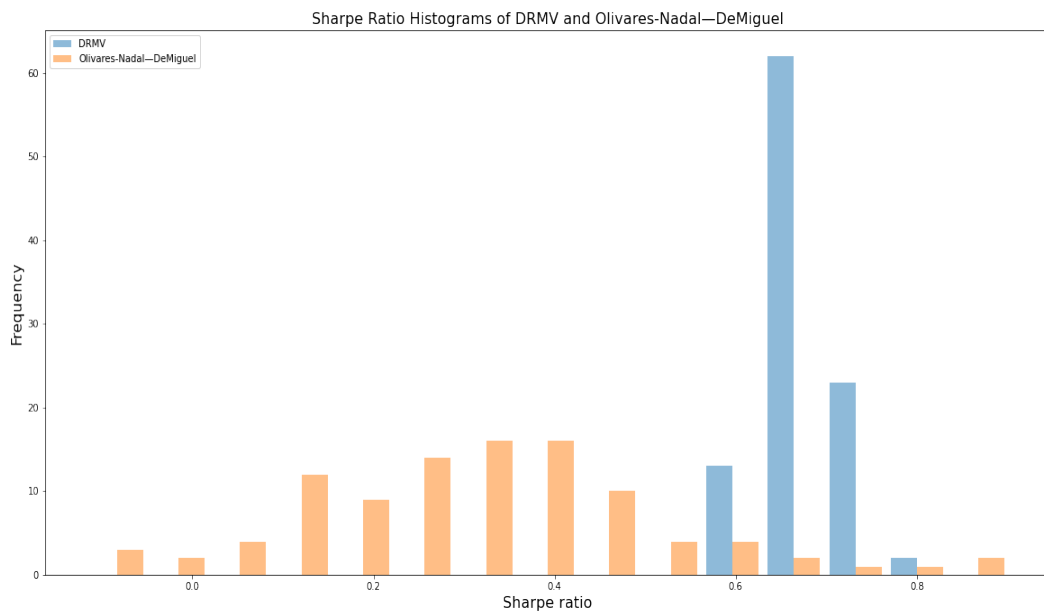


Figure 10: This graph presents the histograms of the Sharpe ratios of the 100 different experiments on DRMV (blue) and Olivares-Nadal-DeMiguel (orange) portfolios. The  $x$ -axis represents the Sharpe ratios and the  $y$ -axis represents the numbers of Sharpe ratios.

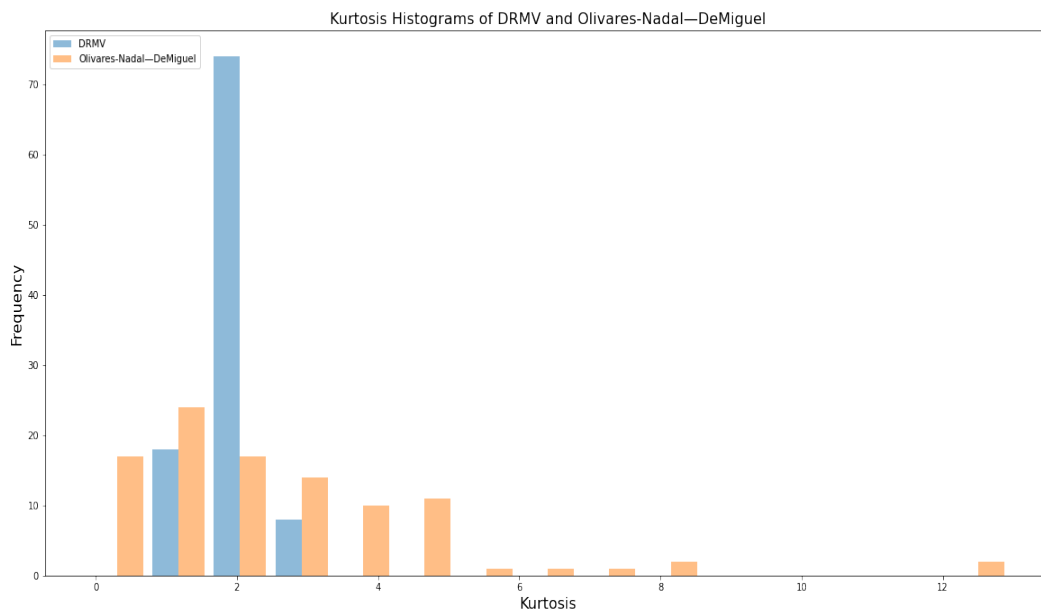


Figure 11: This graph presents the histograms of the kurtosises of the 100 different experiments on DRMV (blue) and Olivares-Nadal-DeMiguel (orange) portfolios. The  $x$ -axis represents the kurtosis and the  $y$ -axis represents the numbers of kurtosises.

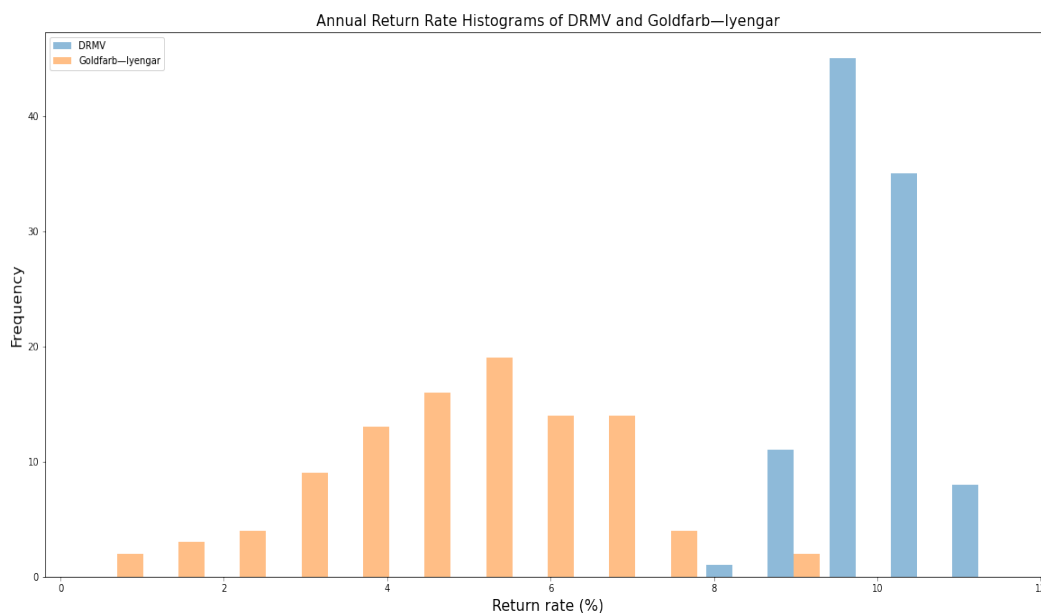


Figure 12: This graph presents the histograms of the annualized returns of the 100 different experiments on DRMV (blue) and Goldfarb-Iyengar (orange) portfolios. There are 2 experiments in which Goldfarb-Iyengar went bankruptcy, which are not included in this histogram. The  $x$ -axis represents the annualized returns and the  $y$ -axis represents the numbers of returns.

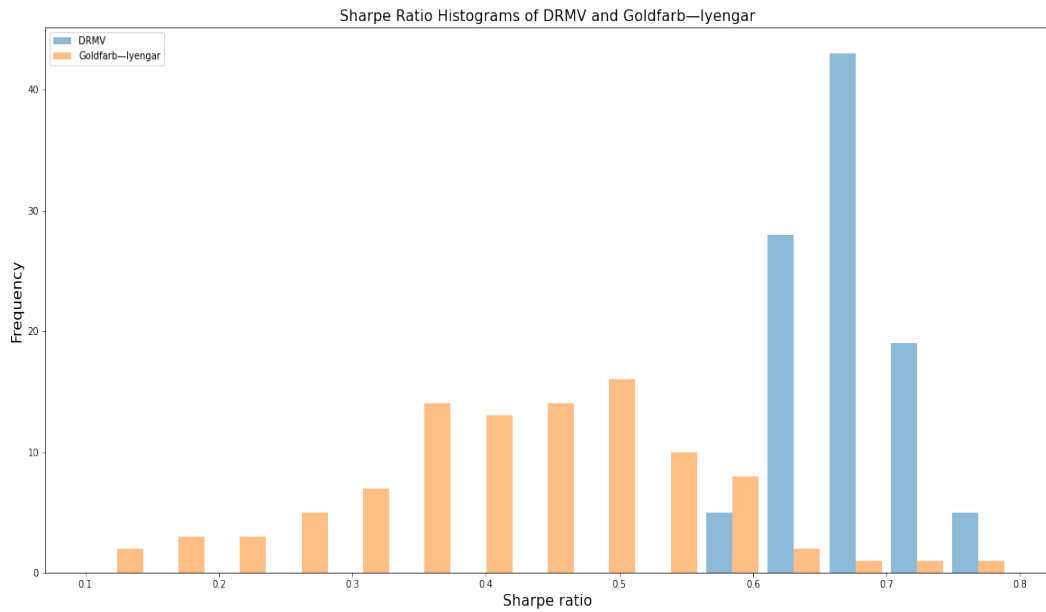


Figure 13: This graph presents the histograms of the Sharpe ratios of the 100 different experiments on DRMV (blue) and Goldfarb-Iyengar (orange) portfolios. The  $x$ -axis represents the Sharpe ratios and the  $y$ -axis represents the numbers of Sharpe ratios.

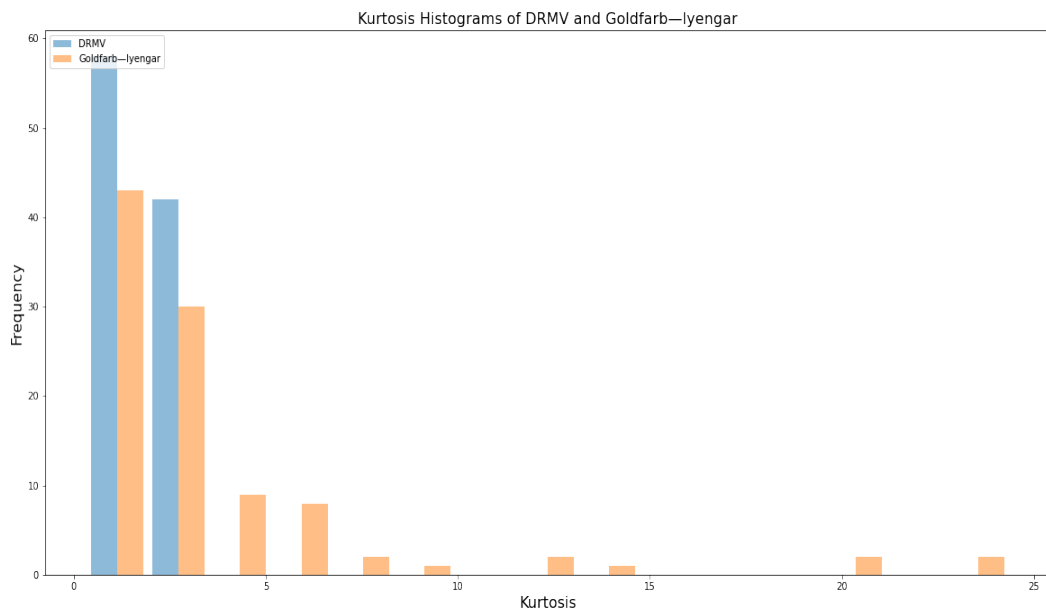


Figure 14: This graph presents the histograms of the kurtosises of the 100 different experiments on DRMV (blue) and Goldfarb-Iyengar (orange) portfolios. The  $x$ -axis represents the kurtosis and the  $y$ -axis represents the numbers of kurtosises.

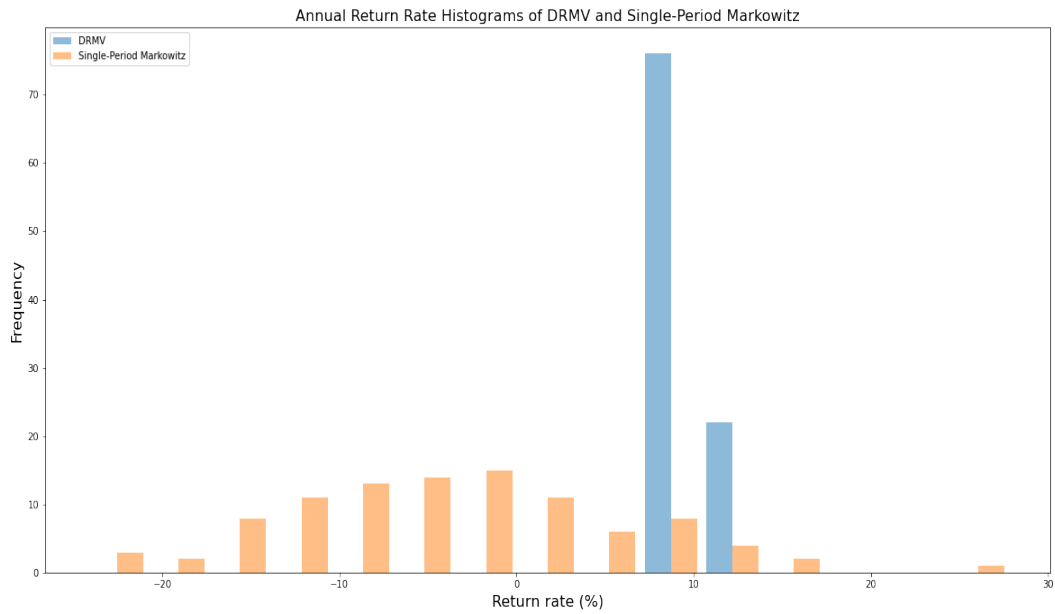


Figure 15: This graph presents the histograms of the annualized returns of the 100 different experiments on DRMV (blue) and single-period Markowitz (orange) portfolios. The  $x$ -axis represents the annualized returns and the  $y$ -axis represents the numbers of returns.

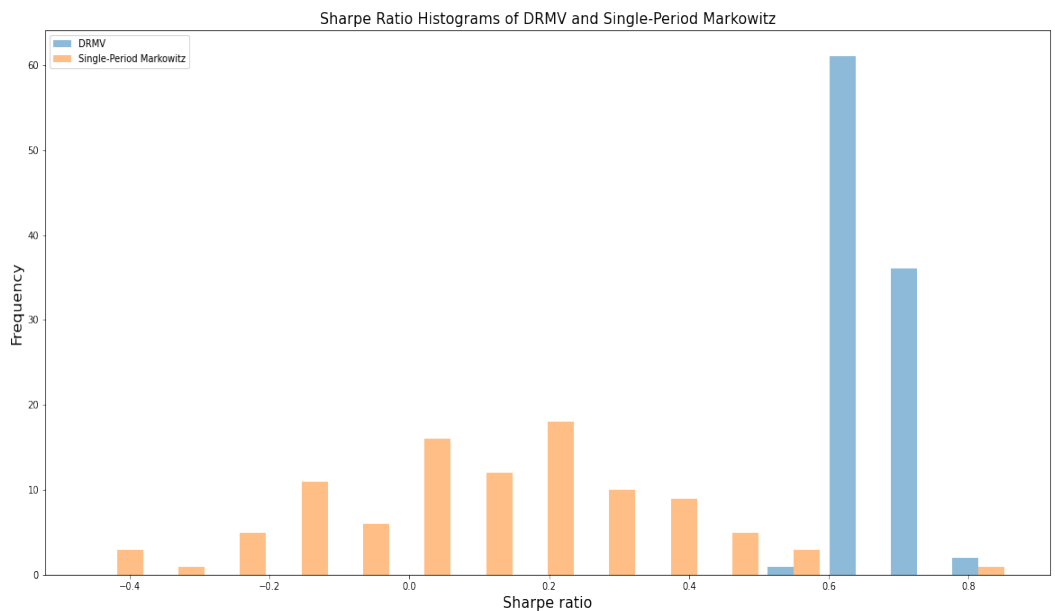


Figure 16: This graph presents the histograms of the Sharpe ratios of the 100 different experiments on DRMV (blue) and single-period Markowitz (orange) portfolios. The  $x$ -axis represents the Sharpe ratios and the  $y$ -axis represents the numbers of Sharpe ratios.

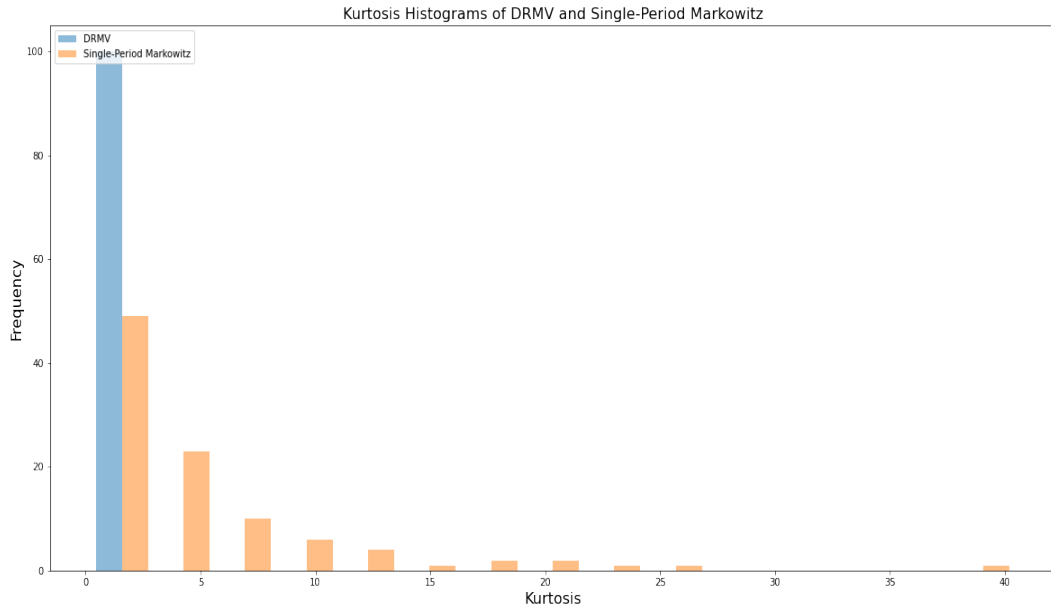


Figure 17: This graph presents the histograms of the kurtosises of the 100 different experiments on DRMV (blue) and single-period Markowitz (orange) portfolios. The  $x$ -axis represents the kurtosis and the  $y$ -axis represents the numbers of kurtosises.

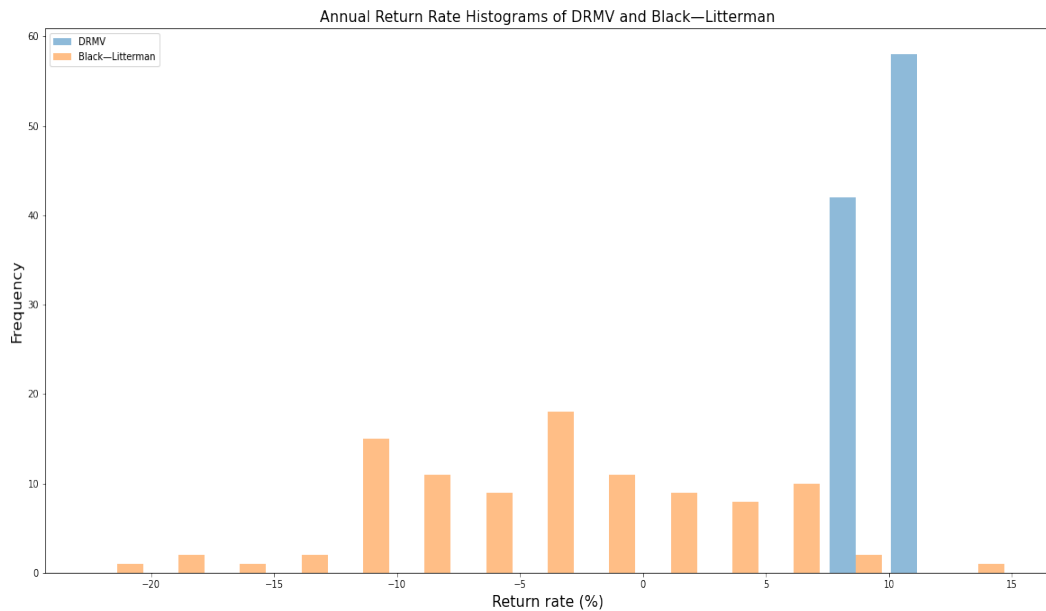


Figure 18: This graph presents the histograms of the annualized returns of the 100 different experiments on DRMV (blue) and Black-Litterman (orange) portfolios. The  $x$ -axis represents the annualized returns and the  $y$ -axis represents the numbers of returns.

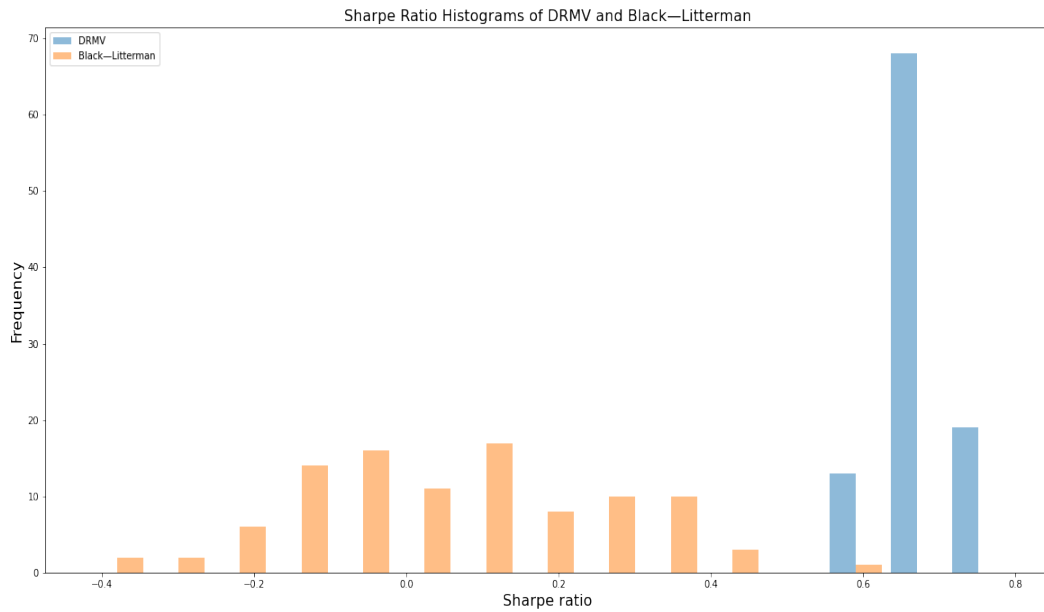


Figure 19: This graph presents the histograms of the Sharpe ratios of the 100 different experiments on DRMV (blue) and Black-Litterman (orange) portfolios. The  $x$ -axis represents the Sharpe ratios and the  $y$ -axis represents the numbers of Sharpe ratios.

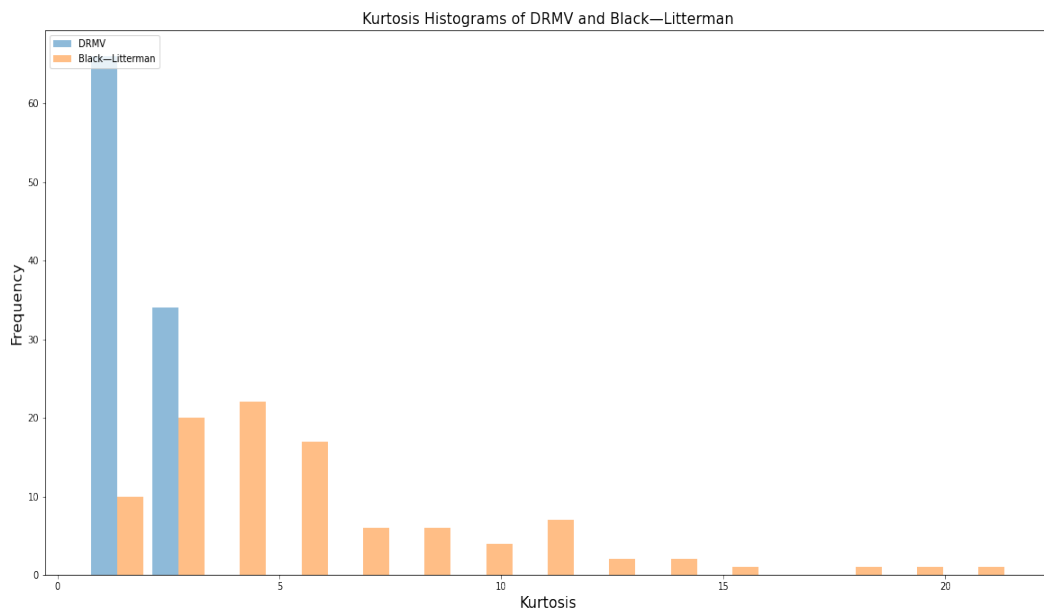


Figure 20: This graph presents the histograms of the kurtosises of the 100 different experiments on DRMV (blue) and Black-Litterman (orange) portfolios. The  $x$ -axis represents the kurtosis and the  $y$ -axis represents the numbers of kurtosises.

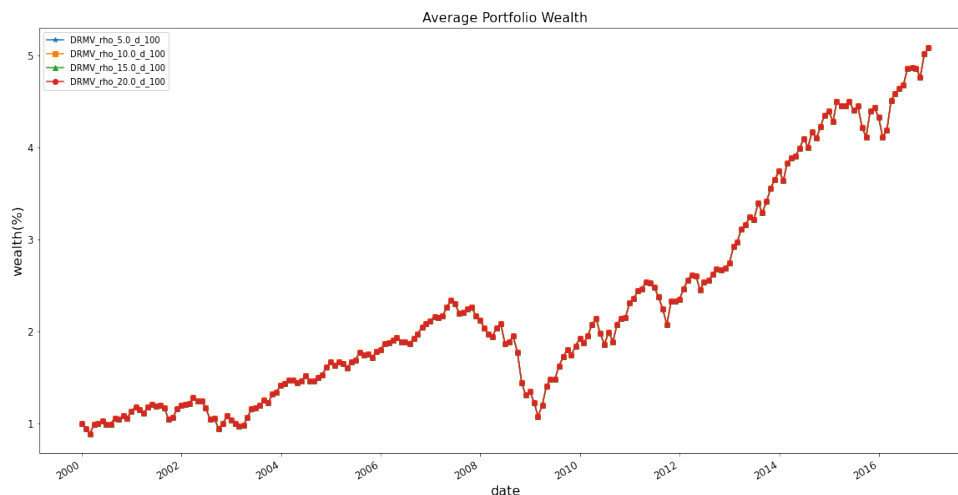


Figure 21: This graph presents DRMV’s average wealth processes from January 2000 to December 2016 with different values of  $\rho$ . The averages are calculated over 100 numerical experiments. The  $x$ -axis indicates the time and the  $y$ -axis indicates the portfolio wealth. Initial wealth is set to be 1.

### 5.3 Discussions

In this subsection we offer discussions on various issues related to our empirical experiments.

#### 5.3.1 Wasserstein order $p = 1$

In all the previously reported experiments we set the order of the Wasserstein distance to be  $p = 2$ . We have also tried  $p = 1$  (and hence  $q = \infty$ ), and found that the performance of the resulting strategies becomes very volatile. So we recommend using  $p = 2$ .

#### 5.3.2 Different targeted returns

We have also tested different (plausible/reasonable) values of the targeted return  $\rho = 5\%, 15\%, 20\%$  in addition to  $\rho = 10\%$ , and found that DRMV maintains the same outperformance with respect to other models and, indeed, some other models become worse under higher targets. Figure 21 plots DRMV’s average wealth processes under these four values of  $\rho$  when  $d = 100$ . They are almost identical so one could probably see only one plot in Figure 21. Hence, the average performance is very robust with different  $\rho$ ’s. This, in turn, suggests that the choice of a specific value of  $\rho$  is unimportant for DRMV so long as it is in the reasonable range of  $[5\%, 20\%]$ , thereby releasing us from the tuning and calibration of this parameter.

#### 5.3.3 Turnover rates

Fabozzi et al. (2007) observe empirically that robust portfolios have low turnover rates. This phenomenon is justified theoretically by the stipulation that robust models are equivalent to



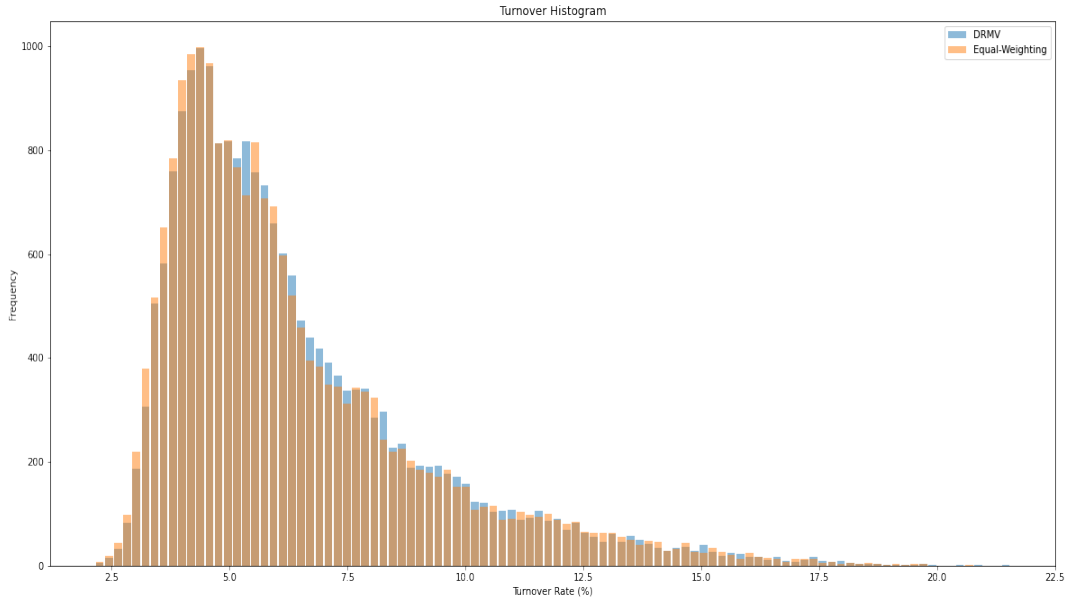


Figure 22: This graph presents the histograms of the DRMV (blue) and equally weighted portfolio (orange) monthly turnover rates of 100 experiments. The  $x$ -axis represents the turnover rate (%) and the  $y$ -axis the numbers of turnover rates.

non-robust models with transaction costs (see the discussion in Section 3) whereas transaction costs in general discourage active trades. We support this with our distributionally robust strategies. Figure 22 shows the histograms of the turnover rates (including buy and sell) with 100 experiments for  $d = 100$  stocks for both DRMV and equally weighted portfolio. We include the latter as it is known to have low turnovers and is a special instance of distributionally robust portfolios.

Clearly, our result reconciles with the finding of Fabozzi et al. (2007). Indeed, the two histograms are almost identical, and most of the monthly turnover rates are lower than 10% which is considered to be very good and reasonable in practice. If we take the average monthly turnover rate to be 10% (definitely an upper bound for the average), then the annual turnover rate is around only 120%.

### 5.3.4 Shrinkage estimators

In all the experiments reported so far, sample covariance was used for the single-period Markovitz model as well as the Black–Litterman model. Upon the recommendation of one of the referees of an earlier version of the paper, we tested using the shrinkage covariance matrices for these two models.

We used the shrinkage estimator of Ollila and Raninen (2018):

$$\Sigma_{\alpha,\beta} = \beta\Sigma_n + \alpha I_n, \quad (21)$$

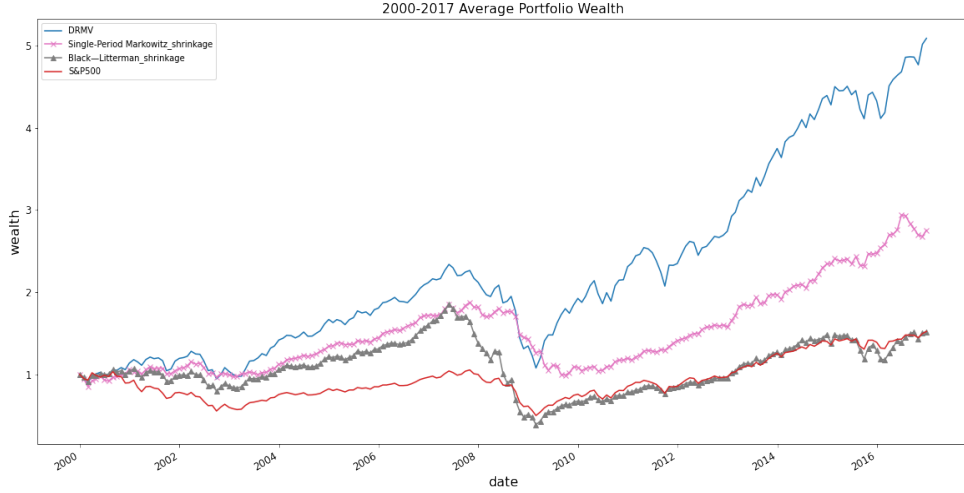


Figure 23: This graph presents portfolios' average wealth processes using shrinkage estimators from January 2000 to December 2016. All the portfolios except S&P 500 consist of 100 stocks and the averages are calculated over 100 numerical experiments. The  $x$ -axis indicates the time in months (from 1 to 204) and the  $y$ -axis indicates the portfolio wealth. Initial wealth is set to be 1.

where  $\Sigma_n$  is the sample covariance matrix and  $I_n$  is the  $n \times n$  identity matrix. The parameters  $\alpha$  and  $\beta$  are estimated by the following

$$\hat{\alpha} = (1 - \hat{\beta})\hat{\eta}, \quad \hat{\beta} = \frac{(\hat{\gamma} - 1)}{(\hat{\gamma} - 1) + \hat{\kappa}(2\hat{\gamma} + d)/n + (\hat{\gamma} + d)/(n - 1)},$$

where  $d$  is the number of stocks,  $\hat{\eta} = \frac{\text{tr}(\Sigma_n)}{d}$ ,  $\hat{\gamma} = \frac{nd}{n-1}[\text{tr}(\Sigma_{sgn}^2) - \frac{1}{n}]$  with  $\Sigma_{sgn} = \frac{1}{n} \sum_{i=1}^n \frac{(R_i - \hat{\mu})(R_i - \hat{\mu})^T}{\|R_i - \hat{\mu}\|}$  and  $\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n \|R_i - \mu\|$ ,  $\hat{\kappa} = \max\left(-\frac{2}{d+2}, \frac{1}{3d} \sum_{j=1}^d \hat{K}_j\right)$  with  $\hat{K}_j = \frac{n-1}{(n-2)(n-3)}[(n+1)\hat{k}_j + 6]$  and  $\hat{k}_j = \frac{m_j^{(4)}}{(m_j^{(2)})^2} - 3$ ,  $m_j^{(q)}$  denoting the  $q$ th order sample moment.

The two models, single-period Markowitz and Black-Litterman, have both been improved with the shrinkage estimators. Figure 23 presents the comparison of average wealth processes between the two models, DRMV and S&P 500. Notably, the Black-Litterman model now outperforms S&P 500 most of the times, as opposed to that without shrinkage (see Figure 2). However, DRMV still dominates all the two models. Histograms of return distributions and Sharpe ratio distributions show superiority of DRMV over the other models similar to those without shrinkage.<sup>19</sup>

<sup>19</sup>As those histograms are similar, we do not present them here.

## 6 Concluding Remarks

We have provided a data-driven distributionally robust theory for Markowitz's mean–variance portfolio selection. The robust model can be solved via a non-robust one based on the empirical probability measure with an additional regularization term. The size of the distributional uncertainty region is not exogenously given; rather it is informed by the return data in a scheme which we have developed in this paper.

Our results may be generalized in different directions. We have chosen the  $l_q$  norm in defining our Wasserstein distance due to its popularity in regularization, but other transportation costs can be used. For example, one may consider the type of transportation cost related to adaptive regularization that has been studied by Blanchet et al. (2017), or the one related to industry cluster as in Blanchet and Kang (2017). Another significant direction is a dynamic (discrete-time or continuous-time) version of the DRMV model.

## Appendices

### A Proof of Theorem 1

The first step is to show that the feasible region over  $\phi$  in the outer minimization part of Problem (2) can be explicitly evaluated. This is given in the following proposition.

**Proposition 2** For  $c(u, v) = \|u - v\|_q^2$ ,  $q \geq 1$ , we have

$$\min_{\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n)} \mathbb{E}_{\mathbb{P}}(\phi^\top R) = \mathbb{E}_{\mathbb{P}_n}(\phi^\top R) - \sqrt{\delta} \|\phi\|_p, \quad (22)$$

where  $p$  satisfies  $1/p + 1/q = 1$ .

**Proof** We consider the following problem

$$\min_{\mathbb{P} \in D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta} \phi^\top \mathbb{E}_{\mathbb{P}}[R] \quad (23)$$

or, equivalently,

$$- \max_{\mathbb{P} \in D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta} \mathbb{E}_{\mathbb{P}}[(-\phi)^\top R]. \quad (24)$$

By checking Slater's condition and using Proposition 4 of Blanchet et al. (2016) we obtain the dual problem:

$$\max_{\mathbb{P} \in D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta} \mathbb{E}_{\mathbb{P}}[(-\phi)^\top R] = \inf_{\lambda \geq 0} \left[ \lambda \delta + \frac{1}{n} \sum_{i=1}^n \Phi_\lambda(R_i) \right] \quad (25)$$

where

$$\begin{aligned} \Phi_\lambda(R_i) &= \sup_u \{h(u) - \lambda c(u, R_i)\} \\ &= \sup_u \{(-\phi^\top)u - \lambda \|u - R_i\|_q^2\} \\ &= \sup_\Delta \{(-\phi^\top)(\Delta + R_i) - \lambda \|\Delta\|_q^2\} \\ &= \sup_\Delta \{(-\phi^\top)\Delta - \lambda \|\Delta\|_q^2\} - \phi^\top R_i \\ &= \sup_\Delta \{\|\phi\|_p \|\Delta\|_q - \lambda \|\Delta\|_q^2\} - \phi^\top R_i \\ &= \frac{\|\phi\|_p^2}{4\lambda} - \phi^\top R_i. \end{aligned}$$

Thus, (25) becomes

$$\begin{aligned} \max_{\mathbb{P} \in D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta} \mathbb{E}_{\mathbb{P}}[(-\phi)^\top R] &= \inf_{\lambda \geq 0} \left\{ \lambda \delta + \frac{1}{n} \sum_{i=1}^n \left[ \frac{\|\phi\|_p^2}{4\lambda} - \phi^\top R_i \right] \right\} \\ &= \inf_{\lambda \geq 0} \left\{ \lambda \delta + \frac{\|\phi\|_p^2}{4\lambda} - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] \right\} \\ &= \sqrt{\delta} \|\phi\|_p - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] \end{aligned}$$

or

$$\min_{\mathbb{P} \in D_c(\mathbb{P}, \mathbb{P}_n) \leq \delta} \phi^\top \mathbb{E}_{\mathbb{P}}[R] = \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] - \sqrt{\delta} \|\phi\|_p. \quad (26)$$

■

Therefore, the feasible region can be rewritten as

$$\mathcal{F}_{\delta, \bar{\alpha}}(n) = \{\phi \in \mathbb{R}^d : \phi^\top \mathbf{1} = 1, \mathbb{E}_{\mathbb{P}_n}(\phi^\top R) \geq \bar{\alpha} + \sqrt{\delta} \|\phi\|_p\},$$

which can now be seen as clearly convex.

Next, by fixing  $\mathbb{E}_{\mathbb{P}}(\phi^\top R) = \alpha \geq \bar{\alpha}$  in the inner maximization part of problem (2) we obtain the following equivalent formulation

$$\min_{\phi \in \mathcal{F}_{\delta, \bar{\alpha}}} \left\{ \max_{\alpha \geq \bar{\alpha}} \left[ \max_{\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n), \mathbb{E}_{\mathbb{P}}(\phi^\top R) = \alpha} \{\phi^\top \mathbb{E}_{\mathbb{P}}(RR^\top) \phi\} - \alpha^2 \right] \right\}. \quad (27)$$

Introducing  $\mathbb{E}_{\mathbb{P}}(\phi^\top R) = \alpha$  is useful because the inner-most maximization problem in the above is now linear in  $\mathbb{P}$ . So, let us concentrate on the problem

$$\max_{\mathbb{P} \in \mathcal{U}_\delta(\mathbb{P}_n), \mathbb{E}_{\mathbb{P}}(\phi^\top R) = \alpha} \phi^\top \mathbb{E}_{\mathbb{P}}(RR^\top) \phi. \quad (28)$$

The following proposition solves this problem in terms of a general cost function  $c$ .

**Proposition 3** *For any cost function  $c$  that is lower semicontinuous and non-negative, the optimal value function of problem (28) is given by*

$$\inf_{\lambda_1 \geq 0, \lambda_2} \left[ \frac{1}{n} \sum_{i=1}^n \Phi(R_i) + \lambda_1 \delta + \lambda_2 \alpha \right], \quad (29)$$

where

$$\Phi(R_i) := \sup_u \left[ (\phi^\top u)^2 - \lambda_1 c(u, R_i) - \lambda_2 \phi^\top u \right].$$

**Proof** The proof is based on a duality argument. Introducing a slack random variable  $S \equiv v$ , where  $v$  is a deterministic number. Then we can recast problem (28) as

$$\max \{ \mathbb{E}_{\mathbb{P}}[(U^\top \phi)^2] : \mathbb{E}_{\pi}[c(U, R) + S] = \delta, \pi_R = \mathbb{P}_n, \pi(S = v) = 1, \quad (30)$$

$$\mathbb{E}_{\pi}[U^\top \phi] = \alpha, \pi \in \mathcal{P}(\mathcal{R}^m \times \mathcal{R}^m \times \mathcal{R}_+) \}. \quad (31)$$

Define

$$\Omega := \{(u, r, s) : c(u, r) < \infty, s \geq 0, r \in \{R_1, \dots, R_n\}\},$$

and let

$$f(u, r, s) = \begin{bmatrix} 1_{r=R_1}(u, r, s) \\ \dots \\ 1_{r=R_n}(u, r, s) \\ \phi^\top u \\ 1_{s=v}(u, r, s) \\ c(u, r) + s \end{bmatrix} \quad \text{and} \quad q = \begin{bmatrix} \frac{1}{n} \\ \dots \\ \frac{1}{n} \\ \alpha \\ 1 \\ \delta \end{bmatrix}. \quad (32)$$

Thus (30) can be written as,

$$\max\{\mathbb{E}_\pi[(U^\top \phi)^2] : \mathbb{E}_\pi[f(U, R, S)] = q, \pi \in \mathcal{P}_\Omega\}. \quad (33)$$

Let  $f_0 = \mathbf{1}_\Omega$ ,  $\tilde{f} = (f_0, f)$ ,  $\tilde{q} = (1, q)$ ,  $\mathcal{Q}_{\tilde{f}} := \{\int \tilde{f}(x)d\mu(x) : \mu \in \mathcal{M}_\Omega^+\}$  where  $\mathcal{M}_\Omega^+$  denote the set of non-negative measures on  $\Omega$ . If  $\phi \neq 0$ , then it is easy to see that  $\tilde{q}$  lies in the interior of  $\mathcal{Q}_{\tilde{f}}$ . By Proposition 6 in Blanchet et al. (2016), the optimal value of problem (33) equals to that of its dual problem, i.e.,

$$\begin{aligned} & \max\{\mathbb{E}_\pi[(U^\top \phi)^2] : \mathbb{E}_\pi[f(U, R, S)] = q, \pi \in \mathcal{P}_\Omega\} \\ &= \inf_{a=(a_0, \dots, a_{n+3}) \in A} \{a_0 + \frac{1}{n} \sum_{i=1}^n a_i + \alpha a_{n+1} + a_{n+2} + \delta a_{n+3}\}, \end{aligned} \quad (34)$$

where

$$\begin{aligned} A := \{a = (a_0, \dots, a_{n+3}) \in \mathbb{R}^{n+4} : a_0 + \frac{1}{n} \sum_{i=1}^n a_i 1_{r=R_i}(u, r, s) + a_{n+1} \phi^\top u \\ + a_{n+2} 1_{s=v}(u, r, s) + a_{n+3} [c(u, r) + s] \geq (\phi^\top u)^2, \forall (u, r, s) \in \Omega\}. \end{aligned}$$

From the definition of  $A$ , replacing  $r = R_i$ , we obtain that the inequality

$$a_0 + a_i + a_{n+2} \geq \sup_{(u, s) \in \Omega} \{(\phi^\top u)^2 - a_{n+3} [c(u, R_i) + s] - a_{n+1} \phi^\top u\} \quad (35)$$

holds for each  $i \in \{1, \dots, n\}$ . It follows directly that

$$\sup_{(u, s) \in \Omega} \{(\phi^\top u)^2 - a_{n+3} [c(u, R_i) + s] - a_{n+1} \phi^\top u\} \quad (36)$$

$$= \begin{cases} +\infty, & \text{if } a_{n+3} < 0 \\ \sup_u \{(\phi^\top u)^2 - a_{n+3} c(u, R_i) - a_{n+1} \phi^\top u\}, & \text{if } a_{n+3} \geq 0. \end{cases} \quad (37)$$

Thus, the dual problem can be expressed as

$$\inf\{a_0 + \frac{1}{n} \sum_{i=1}^n a_i + \alpha a_{n+1} + a_{n+2} + \delta a_{n+3} : a_{n+3} \geq 0, a_0 + a_i + a_{n+2} \geq \sup_u \{(\phi^\top u)^2 - a_{n+3} c(u, R_i) - a_{n+1} \phi^\top u\}\} \quad (38)$$

which can be transformed into

$$\inf_{a_{n+3} \geq 0} \left\{ \frac{1}{n} \sum_{i=1}^n \Phi(R_i) + \alpha a_{n+1} + \delta a_{n+3} \right\}, \quad (39)$$

with

$$\Phi(R_i) := \sup_u \{ (\phi^\top u)^2 - a_{n+3} c(u, R_i) - a_{n+1} \phi^\top u \}.$$

Using  $\lambda_1$  to replace  $a_{n+3}$  and  $\lambda_2$  to replace  $a_{n+1}$ , the dual problem becomes

$$\inf_{\lambda_1 \geq 0} \left\{ \frac{1}{n} \sum_{i=1}^n \Phi(R_i) + \lambda_2 \alpha + \lambda_1 \delta \right\} \quad (40)$$

where

$$\Phi(R_i) := \sup_u \{ (\phi^\top u)^2 - \lambda_1 c(u, R_i) - \lambda_2 \phi^\top u \}.$$

■

Thanks to this proposition, we are able to reduce the inner (infinite dimensional) optimization problem in (2) into a two-dimensional optimization problem in terms of  $\lambda_1$  and  $\lambda_2$ , which can be further simplified if the cost function  $c$  has additional structure. We make this statement precise in the case of a quadratic  $l_q$  cost.

**Proposition 4** *Let  $c(u, v) = \|u - v\|_q^2$  with  $q \geq 1$  and  $1/p + 1/q = 1$ . If  $(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2 - \delta \|\phi\|_p^2 \leq 0$ , then the value of (28) is equal to*

$$\begin{aligned} h(\alpha, \phi) := & \mathbb{E}_{\mathbb{P}_n} \left[ (\phi^\top R)^2 \right] + 2(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R]) \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] + \delta \|\phi\|_p^2 \\ & + 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi}. \end{aligned}$$

**Proof** Writing  $\Delta := u - R_i$ , we have

$$\begin{aligned} \Phi(R_i) &= \sup_u \{ (\phi^\top u)^2 - \lambda_1 c(u, R_i) - \lambda_2 \phi^\top u \} \\ &= \sup_u \{ (\phi^\top u)^2 - \lambda_1 \|u - R_i\|_q^2 - \lambda_2 \phi^\top u \} \\ &= \sup_{\Delta} \{ (\phi^\top (\Delta + R_i))^2 - \lambda_1 \|\Delta\|_q^2 - \lambda_2 \phi^\top (R_i + \Delta) \} \\ &= \sup_{\Delta} \{ (\phi^\top R_i)^2 + (\phi^\top \Delta)^2 + 2(\phi^\top R_i)(\phi^\top \Delta) - \lambda_1 \|\Delta\|_q^2 - \lambda_2 \phi^\top (R_i + \Delta) \} \\ &= (\phi^\top R_i)^2 - \lambda_2 \phi^\top R_i + \sup_{\Delta} \{ (\phi^\top \Delta)^2 + 2(\phi^\top R_i)(\phi^\top \Delta) - \lambda_1 \|\Delta\|_q^2 - \lambda_2 \phi^\top \Delta \} \\ &= (\phi^\top R_i)^2 - \lambda_2 \phi^\top R_i + \sup_{\Delta} \{ (\|\phi\|_p^2 - \lambda_1) \|\Delta\|_q^2 + 2(R_i^\top \phi) - \lambda_2 (\|\phi\|_p \|\Delta\|_q) \}. \end{aligned}$$

We can consider four cases: 1)  $\|\phi\|_p^2 > \lambda_1$ ,  $\Phi(R_i) = +\infty$ ; 2)  $\|\phi\|_p^2 = \lambda_1$ ,  $2R_i^\top \phi \neq \lambda_2$ ,  $\Phi(R_i) = +\infty$ ; 3)  $\|\phi\|_p^2 = \lambda_1$ ,  $2R_i^\top \phi = \lambda_2$ ,  $\Phi(R_i) = 0$ ; 4)  $\|\phi\|_p^2 < \lambda_1$ ,  $\Phi(R_i) = (\phi^\top R_i)^2 - \lambda_2 \phi^\top R_i + \frac{(2R_i^\top \phi - \lambda_2)^2 \|\phi\|_p^2}{4(\lambda_1 - \|\phi\|_p^2)}$ .

For any of the first three cases, the value of  $\frac{1}{n} \sum_{i=1}^n \Phi(R_i)$  is  $+\infty$ . Hence only the fourth case is non-trivial. In this case, problem (29) is transformed into

$$\begin{aligned} & \inf_{\lambda_1 \geq 0, \lambda_2} \left[ \frac{1}{n} \sum_{i=1}^n \Phi(R_i) + \lambda_2 \alpha + \lambda_1 \delta \right] \\ &= \inf_{\lambda_1 \geq \|\phi\|_p^2, \lambda_2} \left\{ \frac{1}{n} \sum_{i=1}^n \left[ (\phi^\top R_i)^2 - \lambda_2 \phi^\top R_i + \frac{(2R_i^\top \phi - \lambda_2)^2 \|\phi\|_p^2}{4(\lambda_1 - \|\phi\|_p^2)} \right] + \lambda_2 \alpha + \lambda_1 \delta \right\} \end{aligned} \quad (41)$$

Define

$$H = \frac{1}{n} \sum_{i=1}^n \left[ (\phi^\top R_i)^2 - \lambda_2 \phi^\top R_i + \frac{(2R_i^\top \phi - \lambda_2)^2 \|\phi\|_p^2}{4(\lambda_1 - \|\phi\|_p^2)} \right] + \lambda_2 \alpha + \lambda_1 \delta.$$

Taking partial derivative with respect to  $\lambda_2$  and setting it to be 0, we get

$$\frac{\partial H}{\partial \lambda_2} = \alpha - \frac{1}{n} \sum_{i=1}^n \left[ \phi^\top R_i + \frac{(2\phi^\top R_i - \lambda_2) \|\phi\|_p^2}{2(\lambda_1 - \|\phi\|_p^2)} \right] = 0$$

which implies (note that  $\phi^\top 1 = 1$  guarantees that  $\|\phi\|_p^2 > 0$ )

$$\lambda_2 = 2\alpha - 2C \frac{\lambda_1}{\|\phi\|_p^2} \quad (42)$$

where  $C := \alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R]$ . Moreover,  $\lambda_2$  is optimal because

$$\frac{\partial^2 H}{\partial \lambda_2^2} = \frac{\|\phi\|_p^2}{2(\lambda_1 - \|\phi\|_p^2)} > 0. \quad (43)$$

We plug (42) into (41) and obtain

$$\begin{aligned} & \inf_{\lambda_1 \geq 0, \lambda_2} \left[ \frac{1}{n} \sum_{i=1}^n \Phi(R_i) + \lambda_2 \alpha + \lambda_1 \delta \right] \\ &= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + \inf_{\lambda_1 \geq \|\phi\|_p^2, \lambda_2} \left\{ \frac{1}{n} \sum_{i=1}^n \left[ -\lambda_2 \phi^\top R_i + \frac{(2R_i^\top \phi - \lambda_2)^2 \|\phi\|_p^2}{4(\lambda_1 - \|\phi\|_p^2)} \right] + \lambda_2 \alpha + \lambda_1 \delta \right\} \\ &= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + \inf_{\lambda_1 \geq \|\phi\|_p^2, \lambda_2} \left\{ \frac{1}{n} \sum_{i=1}^n \left[ \frac{(2R_i^\top \phi - \lambda_2)^2 \|\phi\|_p^2}{4(\lambda_1 - \|\phi\|_p^2)} \right] + \lambda_2 C + \lambda_1 \delta \right\} \\ &= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + \inf_{\lambda_1 \geq \|\phi\|_p^2} \left\{ \frac{1}{n} \sum_{i=1}^n \left[ \frac{(2R_i^\top \phi - 2\alpha + 2C \frac{\lambda_1}{\|\phi\|_p^2})^2 \|\phi\|_p^2}{4(\lambda_1 - \|\phi\|_p^2)} \right] + (2\alpha - 2C \frac{\lambda_1}{\|\phi\|_p^2})C + \lambda_1 \delta \right\}. \end{aligned}$$



Writing  $\lambda_1 = \kappa + \|\phi\|_p^2$ , we have

$$\begin{aligned}
& \inf_{\lambda_1 \geq 0, \lambda_2} \left[ \frac{1}{n} \sum_{i=1}^n \Phi(R_i) + \lambda_2 \alpha + \lambda_1 \delta \right] \\
&= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + \inf_{\kappa \geq 0} \left\{ \frac{1}{n} \sum_{i=1}^n \left[ \frac{(R_i^\top \phi - \alpha + C \frac{\kappa + \|\phi\|_p^2}{\|\phi\|_p^2})^2 N}{\kappa} \right] + (2\alpha - 2C \frac{\kappa + \|\phi\|_p^2}{\|\phi\|_p^2}) C + (\kappa + \|\phi\|_p^2) \delta \right\} \\
&= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + \inf_{\kappa \geq 0} \left\{ \frac{C_1^2}{\|\phi\|_p^2} k + 2\|\phi\|_p^2 (\phi^\top \bar{W} - \alpha + C) + \frac{1}{n} \sum_{i=1}^n \frac{(R_i^\top \phi - \alpha + C)^2 \|\phi\|_p^2}{\kappa} \right. \\
&\quad \left. + 2\alpha C - 2C^2 + \kappa(\delta - \frac{2C^2}{\|\phi\|_p^2}) + \|\phi\|_p^2 \delta \right\} \\
&= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + 2\alpha C - 2C^2 + \|\phi\|_p^2 \delta + \inf_{\kappa \geq 0} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{(R_i^\top \phi - \mathbb{E}_{\mathbb{P}_n}[R] \phi)^2 \|\phi\|_p^2}{\kappa} + \kappa(\delta - \frac{C^2}{\|\phi\|_p^2}) \right\}.
\end{aligned}$$

If  $\delta - C^2/\|\phi\|_p^2 < 0$ , then the optimal value of the above problem is  $-\infty$ , which means that the primal problem (28) is not feasible. If  $\delta - C^2/\|\phi\|_p^2 \geq 0$ , then

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + 2\alpha C - 2C^2 + \|\phi\|_p^2 \delta + \inf_{\kappa \geq 0} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{(R_i^\top \phi - \mathbb{E}_{\mathbb{P}_n}[R] \phi)^2 \|\phi\|_p^2}{\kappa} + \kappa(\delta - \frac{C^2}{\|\phi\|_p^2}) \right\} \\
&= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + 2(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R]) \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] + \delta \|\phi\|_p^2 \\
&\quad + 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\frac{1}{n} \phi^\top \sum_{i=1}^n (R_i - \mathbb{E}_{\mathbb{P}_n}[R]) (R_i - \mathbb{E}_{\mathbb{P}_n}[R])^\top \phi} \\
&= \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + 2(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R]) \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] + \delta \|\phi\|_p^2 \\
&\quad + 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}[R] \phi}.
\end{aligned}$$

Thus, problem (28) can be written as

$$\begin{aligned}
& \min_{\phi} \frac{1}{n} \sum_{i=1}^n (\phi^\top R_i)^2 + 2(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R]) \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] + \delta \|\phi\|_p^2 \\
&\quad + 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}[R] \phi},
\end{aligned}$$

subject to  $1^\top \phi = 1$  and  $(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2 - \delta \|\phi\|_p^2 \leq 0$ .  $\blacksquare$

The condition  $(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2 - \delta \|\phi\|_p^2 \leq 0$  is to make sure that (28) is feasible, failing which the optimal value  $h(\alpha, \phi) = -\infty$ . Proposition 4 ultimately leads to the following main result of the paper, one that transforms (2) into a non-robust portfolio selection problem in terms of the empirical measure  $\mathbb{P}_n$ , with an additional ‘‘regularization’’ term.

We are now ready to prove Theorem 1.

**Proof of Theorem 1.** Note that

$$\begin{aligned}
& h(\alpha, \phi) - \alpha^2 \\
&= \mathbb{E}_{\mathbb{P}_n} \left[ (\phi^\top R)^2 \right] + 2(\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])\phi^\top \mathbb{E}_{\mathbb{P}_n}[R] - \alpha^2 + \delta \|\phi\|_p^2 \\
&+ 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} \\
&= \mathbb{E}_{\mathbb{P}_n} \left[ (\phi^\top R)^2 \right] + 2\alpha \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] - (\phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2 - \alpha^2 - (\phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2 + \delta \|\phi\|_p^2 \\
&+ 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} \\
&= \phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi + \{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2\} \\
&+ 2\sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} \\
&= \left( \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} + \sqrt{\delta \|\phi\|_p^2 - (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2} \right)^2.
\end{aligned}$$

Therefore, it follows from Proposition 4 that

$$\max_{\alpha \geq \bar{\alpha}, (\alpha - \phi^\top \mathbb{E}_{\mathbb{P}_n}[R])^2 - \delta \|\phi\|_p^2 \leq 0} [h(\alpha, \phi) - \alpha^2] = \left( \sqrt{\phi^\top \text{Var}_{\mathbb{P}_n}(R) \phi} + \sqrt{\delta \|\phi\|_p^2} \right)^2,$$

with the optimal  $\alpha_{opt} = \phi^\top \mathbb{E}_{\mathbb{P}_n}[R] \geq \bar{\alpha}$ . This concludes the proof.

## B Proof of Theorem 2

Define

$$h_0(R, \Sigma) = RR^\top - \Sigma \quad \text{and} \quad h_1(R, \mu) = R - \mu.$$

Then, by Proposition 1 of Blanchet et al. (2016) we have that for any given  $\mu$  and  $\Sigma$ ,

$$\mathcal{R}_n(\Sigma, \mu) = \sup_{\Lambda \in \mathbb{R}^{d \times d}, \lambda \in \mathbb{R}^d} \left\{ -\mathbb{E}_{\mathbb{P}_n} \left[ \sup_{u \in \mathbb{R}^d} \{ \text{Tr}(\Lambda h_0(u, \Sigma)) + \lambda^\top h_1(u, \mu) - \|u - R\|_q^2 \} \right] \right\}.$$

Observe that

$$\begin{aligned}
& \sup_{u \in \mathbb{R}^d} \{ \text{Tr}(\Lambda h_0(u, \Sigma)) + \lambda^\top h_1(u, \mu) - \|u - R\|_q^2 \} \\
&= \sup_{\Delta \in \mathbb{R}^d} \{ \text{Tr}(\Lambda h_0(\Delta + R, \Sigma)) + \lambda^\top h_1(\Delta + R, \mu) - \|\Delta\|_q^2 \} \\
&= \sup_{\Delta \in \mathbb{R}^d} \{ \text{Tr}(\Lambda [h_0(\Delta + R, \Sigma) - h_0(R, \Sigma)]) + \lambda^\top \Delta - \|\Delta\|_q^2 \} \\
&+ \text{Tr}(\Lambda h_0(R, \Sigma)) + \lambda^\top h_1(R, \mu).
\end{aligned}$$

Moreover, let us write

$$\text{Tr}(\Lambda [h_0(\Delta + R, \Sigma) - h_0(R, \Sigma)]) = \int_0^1 \frac{d}{dt} \text{Tr}(\Lambda h_0(R + t\Delta)) dt.$$

However,

$$\begin{aligned}\frac{d}{dt}Tr(\Lambda h_0(R+t\Delta)) &= 2Tr(\Lambda(R+t\Delta)\Delta^\top) \\ &= 2Tr(\Lambda R\Delta^\top) + 2t\Delta^\top\Lambda\Delta.\end{aligned}$$

Furthermore,

$$\mathbb{E}_{\mathbb{P}_n}[Tr(\Lambda h_0(R, \Sigma))]|_{\Sigma=\Sigma_n} = 0. \quad (44)$$

So, we deduce

$$\begin{aligned}\mathcal{R}_n(\Sigma_n, \mu) &= \sup_{\lambda \in \mathbb{R}^d} \{-\mathbb{E}_{\mathbb{P}_n}[\lambda^\top(R-\mu)] + \\ &\quad \sup_{\Lambda \in \mathbb{R}^{d \times d}} (-\mathbb{E}_{\mathbb{P}_n}[\sup_{\Delta} \{2Tr(\Lambda R\Delta^\top) + \Delta^\top\Lambda\Delta + \lambda^\top\Delta - \|\Delta\|_q^2\}])\}.\end{aligned}$$

Introduce the scaling  $\Delta = \bar{\Delta}/n^{1/2}$  and  $\bar{\lambda} = \lambda n^{1/2}$  and  $\bar{\Lambda} = \Lambda n^{1/2}$ . Then we obtain

$$\begin{aligned}n\mathcal{R}_n(\Sigma_n, \mu_n) &= \sup_{\bar{\lambda} \in \mathbb{R}^d} \{-n^{-1/2} \sum_{i=1}^n \bar{\lambda}^\top(R_i - \mu_n) + \\ &\quad \sup_{\bar{\Lambda} \in \mathbb{R}^{d \times d}} (-\mathbb{E}_{\mathbb{P}_n}[\sup_{\bar{\Delta}} \{2Tr(\bar{\Lambda} R\bar{\Delta}^\top) + \bar{\Delta}^\top\bar{\Lambda}\bar{\Delta}/n^{1/2} + \bar{\lambda}^\top\bar{\Delta} - \|\bar{\Delta}\|_q^2\}])\}.\end{aligned}$$

In the proof of Proposition 3 in Blanchet et al. (2016), under Assumption A2), a technique is introduced to show that  $\bar{\Delta}$  and  $\bar{\lambda}$  can be restricted to compact sets with high probability and therefore the term  $\bar{\Delta}^\top\bar{\Lambda}\bar{\Delta}/n^{1/2}$  is asymptotically negligible. On the other hand,

$$\begin{aligned}\sup_{\bar{\Delta}} \{2Tr(\bar{\Delta}^\top\bar{\Lambda}R) + \bar{\Delta}^\top\bar{\lambda} - \|\bar{\Delta}\|_q^2\} \\ = \sup_{\bar{\Delta}} \{2\|\bar{\Lambda}R + \bar{\lambda}\|_p \|\bar{\Delta}\|_q - \|\bar{\Delta}\|_q^2\} = \|\bar{\Lambda}R + \bar{\lambda}\|_p^2.\end{aligned}$$

Therefore, if

$$n^{-1/2} \sum_{i=1}^n (R_i - \mu_n) \Rightarrow -Z$$

for some  $Z$  (to be characterized momentarily), then we conclude that

$$\mathcal{R}_n(\Sigma_n, \mu_n) \Rightarrow L_0 = \sup_{\bar{\lambda} \in \mathbb{R}^d} \{\bar{\lambda}^\top Z - \inf_{\bar{\Lambda} \in \mathbb{R}^{d \times d}} \mathbb{E}_{\mathbb{P}^*}[\|\bar{\Lambda}R + \bar{\lambda}\|_p^2]\}.$$

If  $p = 2$  then we have

$$\mathbb{E}_{\mathbb{P}^*}[\|\bar{\Lambda}R + \bar{\lambda}\|_2^2] = \sum_i \mathbb{E}_{\mathbb{P}^*}(\bar{\Lambda}_i \cdot R + \bar{\lambda}_i)^2.$$

So, taking derivative with respect to the  $i$ -th row,  $\bar{\Lambda}_i$ , of the matrix  $\bar{\Lambda}$ ,  $\bar{\Lambda}_i$ , we obtain

$$\nabla_{\bar{\Lambda}_i} \mathbb{E}_{\mathbb{P}^*} [\|\bar{\Lambda}R + \bar{\lambda}\|_2^2] = 2\mathbb{E}_{\mathbb{P}^*} \left( \left( R^\top \bar{\Lambda}_i + \bar{\lambda}_i \right) R \right) = 2\mathbb{E}_{\mathbb{P}^*} \left( R^\top \bar{\Lambda}_i R \right) + 2\bar{\lambda}_i \mathbb{E}_{\mathbb{P}^*} (R) = 0. \quad (45)$$

Writing

$$\mu_* = \mathbb{E}_{\mathbb{P}^*} (R) \text{ and } \Sigma_* = \mathbb{E}_{\mathbb{P}^*} (RR^\top),$$

we obtain

$$\Sigma_* \bar{\Lambda}_i = -\bar{\lambda}_i \mu_*.$$

Solving this equation yields  $\bar{\Lambda}_i = -\bar{\lambda}_i \Sigma_*^{-1} \mu_*$ . Therefore,

$$\begin{aligned} \mathbb{E}_{\mathbb{P}^*} (\bar{\Lambda}_i R + \bar{\lambda}_i)^2 &= \bar{\Lambda}_i^\top \Sigma_* \bar{\Lambda}_i + 2\bar{\lambda}_i \bar{\Lambda}_i^\top \mu_* + \bar{\lambda}_i^2 \\ &= \bar{\lambda}_i^2 \left( 1 - \mu_*^\top \Sigma_*^{-1} \mu_* \right). \end{aligned}$$

By assumption A3),  $\text{Var}_{\mathbb{P}^*} (R)$  is positive definite; hence invertible. It then follows from the ShermanMorrison formula that

$$\begin{aligned} (\Sigma_*)^{-1} &= (\text{Var}_{\mathbb{P}^*} (R) + \mu_* \mu_*^\top)^{-1} \\ &= \text{Var}_{\mathbb{P}^*} (R)^{-1} - \frac{\text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_* \mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1}}{1 + \mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_*}. \end{aligned}$$

So

$$\begin{aligned} \mu_*^\top (\Sigma_*)^{-1} \mu_* &= \mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_* - \frac{(\mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_*)^2}{1 + \mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_*} \\ &= \frac{\mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_*}{1 + \mu_*^\top \text{Var}_{\mathbb{P}^*} (R)^{-1} \mu_*} < 1, \end{aligned}$$

leading to

$$\begin{aligned} L_0 &= \sup_{\bar{\lambda} \in \mathbb{R}^d} \{ \bar{\lambda}^\top Z - \inf_{\bar{\Lambda} \in \mathbb{R}^{d \times d}} \mathbb{E}_{\mathbb{P}^*} [\|\bar{\Lambda}R + \bar{\lambda}\|_2^2] \} \\ &= \sup_{\bar{\lambda} \in \mathbb{R}^d} \{ \bar{\lambda}^\top Z - \|\bar{\lambda}\|_2^2 \left( 1 - \mu_*^\top \Sigma_*^{-1} \mu_* \right) \} \\ &= \frac{\|Z\|_2^2}{4 \left( 1 - \mu_*^\top \Sigma_*^{-1} \mu_* \right)}. \end{aligned}$$

It remains to identify  $Z$ . Observe that

$$\begin{aligned} \mu_n &= \rho 1 + 2 \left( \Sigma_n \phi^* - \phi^{*T} \Sigma_n \phi^* 1 \right) / \lambda_1^* \\ &= \rho 1 + 2 \left( \Sigma_* \phi^* - \phi^{*T} \Sigma_* \phi^* 1 \right) / \lambda_1^* \\ &\quad + 2 \left( H_n \phi^* - \phi^{*T} H_n \phi^* 1 \right) / \lambda_1^* \\ &= \mu_* + 2 \left( H_n \phi^* - \phi^{*T} H_n \phi^* 1 \right) / \lambda_1^*, \end{aligned}$$

where  $H_n := \Sigma_n - \Sigma_*$ . By A1) we have

$$\begin{aligned} n^{-1/2} \sum_{i=1}^n (R_i - \mu_*) &\Rightarrow Z_0 \sim N(0, \Upsilon_{g_1}), \\ n^{1/2} H_n &\Rightarrow Y \sim N(0, \Upsilon_{g_2}). \end{aligned}$$

Thus,

$$\begin{aligned} n^{-1/2} \sum_{i=1}^n \bar{\lambda}^\top (R_i - \mu_*) + 2n^{1/2} \bar{\lambda}^\top (H_n \phi^* - \phi^{*T} H_n \phi^* \mathbf{1}) / \lambda_1^* \\ \Rightarrow \bar{\lambda}^\top Z = \bar{\lambda}^\top (Z_0 + Z_1), \end{aligned}$$

where

$$Z_1 := 2(Y \phi^* - \phi^{*T} Y \phi^* \mathbf{1}) / \lambda_1^*.$$

## C Proof of Proposition 1

Note that

$$\begin{aligned} n^{1/2} \left\{ \frac{(\phi_n^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi_n^*\|_p} \right\} &= n^{1/2} \left\{ \frac{(\phi_n^* - \phi^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi_n^*\|_p} \right\} \\ &\quad + n^{1/2} \left\{ \frac{(\phi^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi^*\|_p} \cdot \frac{\|\phi^*\|_p}{\|\phi_n^*\|_p} \right\}. \end{aligned}$$

By the standard central limit theorem and the fact that  $\phi_n^* \rightarrow \phi^*$  in probability, we conclude

$$n^{1/2} \left\{ \frac{(\phi_n^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi_n^*\|_p} - \frac{(\phi^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi^*\|_p} \right\} \Rightarrow 0$$

as  $n \rightarrow \infty$ . However, again by the central limit theorem we have

$$n^{1/2} \left\{ \frac{(\phi^*)^\top [\mathbb{E}_{\mathbb{P}_n}(R) - \mathbb{E}_{\mathbb{P}^*}(R)]}{\|\phi^*\|_p} \right\} \Rightarrow N(0, \Upsilon_{\phi^*}),$$

which yields the desired result.

## References

- Ao, M., Li, Y., and Zheng, X. Approaching mean-variance efficiency for large portfolios. *The Review of Financial Studies*, 32:2890–2919, 2018.
- Bertsimas, D., Pachamanova, D., and Sim, M. Robust linear optimization under general norms. *Operations Research Letters*, 32:510–516, 2004.

- Bielecki, T., Jin, H., Pliska, S., and Zhou, X. Continuous-time mean–variance portfolio selection with bankruptcy prohibition. *Mathematical Finance*, 15:213–244, 2005.
- Blanchet, J. and Kang, Y. Distributionally robust groupwise regularization estimator. *arXiv:https://arxiv.org/abs/1705.04241*, 2017.
- Blanchet, J. and Murthy, K. Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 44(2):565–600, 2019.
- Blanchet, J., Kang, Y., and Murthy, K. Robust wasserstein profile inference and applications to machine learning. *arXiv:https://arxiv.org/abs/1610.05627*, 2016.
- Blanchet, J., Kang, Y., Zhang, F., and Murthy, K. Data-driven optimal transport cost selection for distributionally robust optimization. *arXiv:https://arxiv.org/abs/1705.07152*, 2017.
- Blanchet, J., Murthy, K., and Si, N. Confidence regions in wasserstein distributionally robust estimation. *https://arxiv.org/abs/1906.01614*, 2019.
- Costa, O. and Paiva, A. Robust portfolio selection using linear-matrix inequalities. *Journal of Economic Dynamics and Control*, 26:889–909, 2002.
- Cui, X., Gao, J., and Li, D. Continuous-time mean-variance portfolio selection with finite transactions. *Stochastic analysis and applications to finance*, 13:77–98, 2012.
- Delage, E. and Ye, Y. Data-driven distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations Research*, 58:595–612, 2010.
- Delage, E., Kuhn, D., and Wiesemann, W. Decision-making under uncertainty: When can a random decision reduce risk? *Management Science*, 65:32823301, 2019.
- DeMiguel, V., Garlappi, L., and Uppal, R. Optimal versus naive diversification: How inefficient is the  $1/n$  portfolio strategy? *The Review of Financial Studies*, 22:1915–1953, 2009.
- Esfahani, P. and Kuhn, D. Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1-2):115–166, 2018.
- Fabozzi, F., Kolm, P., Pachamanova, D., and Focardi, S. Robust portfolio optimization. *Journal of Portfolio Management*, 33(3):40–48, 2007.
- Fama, E. F. and French, K. R. The cross-section of expected stock returns. *The Journal of Finance*, 47(2):427–465, 1992.
- Fama, E. F. and French, K. R. Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics*, 33:3–56, 1993.
- Fan, J., Liao, Y., and Mincheva, M. High-dimensional covariance matrix estimation in approximate factor models. *The Annals of Statistics*, 39:3320–3356, 2011.

- Gao, R. and Kleywegt, A. Distributionally robust stochastic optimization with wasserstein distance. <https://arxiv.org/abs/1604.02199>, 2016.
- Gao, R., Chen, X., and Kleywegt, A. Wasserstein distributional robustness and regularization in statistical learning. <https://arxiv.org/abs/1712.06050>, 2017.
- Ghaoui, L., Oks, M., and Oustry, F. Worst-case value-at-risk and robust portfolio optimization: A conic programming approach. *Operations Research*, 51:543–556, 2003.
- Goh, J. and Sim, M. Distributionally robust optimization and its tractable approximations. *Operations Research*, 58:595–612, 2010.
- Goldfarb, D. and Iyengar, G. Robust portfolio selection problems. *Mathematics of Operations Research*, 28(1):1–38, 2003.
- Gotoh, J. and Takeda, A. On the role of norm constraints in portfolio selection. *Computational Management Science*, 8:323–353, 2011.
- Halldorsson, B. and Tutuncu, R. An interior-point method for a class of saddle-point problems. *Journal of Optimization Theory and Applications*, 116:559–590, 2003.
- Hansen, L. and Sargent, T. *Robustness*. Princeton University Press, Princeton, N.J., 2008.
- Hu, Z. and Hong, L. Kullback-leibler divergence constrained distributionally robust optimization. *Available at Optimization Online*, 2013.
- Idzorek, T. A step-by-step guide to the black-litterman model. *Technical Report*, 2002.
- Jiang, R. and Guan, Y. Data-drive chance constrained stochastic program. *Mathematical Programming*, 158:291–327, 2016.
- Kantorovich, L. On the transfer of masses (in russian). *Doklady Akademii Nauk*, 37:227–229, 1942.
- Lobo, M. and Boyd, S. The worst-case risk of a portfolio. [https://web.stanford.edu/~boyd/papers/pdf/risk\\_bnd.pdf](https://web.stanford.edu/~boyd/papers/pdf/risk_bnd.pdf), 2000.
- Markowitz, H. Portfolio selection. *Journal of Finance*, 7:77–91, 1952.
- Monge, G. Mmoire sur la thorie des dblais et des remblais. *Histoire de l'Academie Royale des Sciences de Paris*, 1781.
- Olivares-Nadal, A. and DeMiguel, V. Technical note-a robust perspective on transaction costs in portfolio optimization. *Operations Research*, 66(3):733–739, 2018.
- Ollila, E. and Raninen, E. Optimal shrinkage covariance matrix estimation under random sampling from elliptical distribution. <https://arxiv.org/abs/1808.10188>, 2018.
- Petersen, I., James, M., and Dupuis, P. Minimax optimal control of stochastic uncertain systems with relative entropy constraints. *IEEE Transactions on Automatic Control*, 45(3):398–412, 2000.

- Pflug, G. and Wozabal, D. Ambiguity in portfolio selection. *Quantitative Finance*, 7:435–442, 2007.
- Villani, C. *Topics in optimal transportation*, volume 58. Graduate Studies in Mathematics, American Mathematics Society, Providence, RI., 2003.
- Wiesemann, W., Kuhn, D., and Sim, M. Distributionally robust convex optimization. *Operations Research*, 62:1358–1376, 2014.
- Wozabal, D. A framework for optimization under ambiguity. *Annals of Operations Research*, 193:21–47, 2012.
- Zhao, C. and Guan, Y. Data-driven risk-averse stochastic optimization with wasserstein metric. *Operations Research Letters*, 46(2):262–267, 2018.